

**APPLICATION UNDER UNITED STATES PATENT LAWS**

Atty. Dkt. No. 009848-0314964

Invention: ISOLATION AND CLONING OF DNA FROM UNCULTIVATED ORGANISMS

Inventor (s): Achim Quaizer;  
 Torsten Ochsenreiter;  
 Alexander H. Treusch;  
 Arnulf Kletzin;

Christa Schleper;  
 Patrick Lorenz;  
 Jürgen Eck

**CERTIFICATE OF EXPRESS MAILING UNDER 37 C.F.R. §1.10**

I hereby certify that this correspondence (along with any paper referred to as being attached or enclosed) is being mailed via "Express Mail Post Office to Addressee" service of the United States Postal Service (Express Mail Label No. EL 989437908 US) on the date shown below in an envelope addressed to the Commissioner for Patents, P.O. Box 1450, Alexandria, VA 22313-1450.

Date: February 18, 2005

By: 

Sachiko Y. Snedden

**Address communications to the  
correspondence address associated  
 with our Customer No**

**27500**

Pillsbury Winthrop LLP

This is a:

- ☐ Provisional Application
- ☐ Regular Utility Application
- ☐ Continuing Application  
☐ The contents of the parent are incorporated by reference
- ☒ PCT National Phase Application
- ☐ Design Application
- ☐ Reissue Application
- ☐ Plant Application
- ☐ Substitute Specification  
Sub. Spec Filed  
 in App. No. /
- ☐ Marked up Specification re  
Sub. Spec. filed  
 In App. No. /

Pillsbury Winthrop LLP  
 Intellectual Property Group  
 11682 El Camino Real, Suite 200  
 San Diego, CA 92130-2092  
 Attorneys  
 Telephone: (619) 234-5000

**SPECIFICATION**

## **ISOLATION AND CLONING OF DNA FROM UNCULTIVATED ORGANISMS**

The present invention relates to a device for the isolation and/or purification of nucleic acid molecules suitable to bind and/or inactivate inhibitors of the activity of reagents or enzymes used for DNA manipulation and to separate a plurality of nucleic acid molecules with respect to their size. Moreover, the invention relates to a method for the isolation of a nucleic acid molecule comprising applying a sample to the device of the invention wherein said nucleic acid molecule preferably is part of a sample which represents a fraction of the metagenome of a given habitat. Furthermore, the invention relates to a method for the generation of at least one gene library comprising nucleic acid molecules isolated by the method of the invention and to a nucleic acid molecule isolated by the method of the invention and with the device of the present invention.

Several documents are cited throughout the text of this specification. The disclosure content of the documents cited herein (including any manufacture's specifications, instructions, etc.) is herewith incorporated by reference.

Enzymes are highly efficient biological catalysts. As such they are key players in environmentally friendly technical conversion processes of modern sustainable biotechnology.

Enzymes particularly from microbial sources are active ingredients in many processes of the textile, detergent, pulp- and paper, food and feed industries. In addition widespread stereoselective substrate recognition and conversion make enzymes particularly attractive for synthetic organic chemists in need of chiral specificity. A bottleneck in the development of innovative technical processes based on enzymes is the supply with suitable new biocatalysts. Owing to their phylogeny and physiological diversity microorganisms constitute the largest resource of natural genetic and enzymatic diversity. However the largest proportion of microorganisms

evades cultivation under laboratory conditions (Amann et. al. (1995). *Microbiol Rev* 59, 143-69). Classic microbiology relying on cultivation of pure strains to provide homogenous and defined systems for homologous enzyme production and to supply genomic DNA for recombinant expression strategies therefore inevitably fails to access the entire biosynthetic potential harboured in this enormous natural resource. The recent development of strategies to directly isolate and clone genomic DNA from non-cultivated microbial consortia opens up new dimensions of accessible enzymatic diversity (Rondon et. al. (2000). *Appl Environ Microbiol* 66, 2541-7). Fundamental work on the handling of DNA from non-cultivated microorganisms - the so-called metagenome – (Handelsman et. al. (1998) *Chem Biol* 5, R245-9) by Torsvik (Torsvik and Goksoyr (1978) *Soil Biology and Biochemistry* 10, 7-12), (Torsvik (1980) *Soil Biology and Biochemistry* 12, 15-21), Somerville (Somerville et. al. (1989) *Appl Environ Microbiol* 55(3), 548-554) and Schmidt (Schmidt et. al. (1991) *Journal of Bacteriology* 173, 4371-4378) showed that genomic DNA can be directly isolated from complex microbial assortments as present, inter alia, in plankton or soil. This DNA may be digested and cloned into suitable vectors for recombinant maintenance in heterologous hosts to generate screenable gene libraries. Such metagenome libraries were shown to be useful in the identification of novel genes from uncultivated organisms. The discovery of novel enzymes by screening of non-normalised metagenome libraries from planktonic and soil sources has been reported in the literature (Cottrell et. al. (1999) *Appl Environ Microbiol* 65, 2553-7), (Henne et. al. (1999) *Appl Environ Microbiol* 65, 3901-7; Henne et. al. (2000) *Appl Environ Microbiol* 66, 3113-3116); (US-patent No : 5,849,491); (Rondon et. al. (2000) *Appl Environ Microbiol* 66, 2541-7). The list of enzyme activities discovered in this way (lipase, esterase, amylase, nuclease, chitinase, xylanase) is still rather small. Importantly also more complex activities like the production of bioactive secondary metabolites requiring entire gene clusters for expression have been identified in metagenomic libraries (MacNeil et. al. (2001) *J Mol Microbiol Biotechnol* 3, 301-8) (Wang et. al. (2000) *Org Lett* 2, 2401-4) (Brady et. al. (2001) *Org Lett* 3, 1981-4). Secondary metabolites, like polyketides, are often produced by enzyme complexes encoded by assortments of genes covering in excess of 100 kbp of contiguous DNA (Schwecke et. al. (1995) *Proc Natl Acad Sci U S A* 92, 7839-43). The cloning of such

large fragments of environmental DNA is much more challenging than cloning smaller DNA fragments and is substantially facilitated by the current invention. Proprietary technology for the cloning particularly of normalised environmental DNA and the screening of libraries generated thereby is described in US patents US 6,280,926; US 6,054,267; US 6,057,103; US 6,001,574 and PCT applications WO99/45154; WO98/58085; WO99/10539.

DNA directly extracted from microbial consortia in the context of their natural substratum usually is contaminated with substances inhibiting standard enzymatic manipulations that are essentially required for cloning, analysis or amplification of nucleic acids carrying genetic information. In particular the efficiencies of DNA digestion with restriction enzymes (Tsai and Olson (1992) *Appl Environ Microbiol* 58, 2292-5), (Tebbe and Vahjen (1993) *Appl Environ Microbiol* 59, 2657-65), the polymerase chain reaction (PCR) (Zhou et. al. (1996) *Appl Environ Microbiol* 62, 316-22), DNA-DNA hybridisation and bacterial transformation with environmental DNA (Tebbe and Vahjen (1993) loc. cit.) are inversely correlated with natural substrate derived inhibitor concentrations. Besides inorganic inhibitors like heavy metal ions, there are polysaccharides and in particular humic and fulvic acids that act as the single most important sources of above mentioned inhibitions. Humic and fulvic acids are high molecular weight heterocyclic polyphenols mainly of plant origin with an affinity to polynucleotides and strongly protein denaturing properties ((Young et. al. (1993) *Appl Environ Microbiol* 59, 1972-1974); see appended figure 1).

Yet, the efficient removal of such inhibitors is a prerequisite for all enzymatic manipulations required, e.g., for cloning DNA, in particular environmental DNA, into suitable vectors. Several strategies have been pursued. Simple dilution of contaminated DNA to bring inhibitor concentrations below a critical threshold may be sufficient if the subsequent enzymatic manipulation is of suitable power to compensate for the concomitant reduction in target/substrate concentration. Surely such dilution will significantly curtail the efficiency of most subsequent molecular manipulations necessary for cloning following simple mass- action laws. The polymerase chain reaction (PCR), owing to its exponential amplification strategy is powerful enough to generate strong signals even from very low target numbers and

often reducing the amount of input environmental DNA (and inhibitors) in a reaction will substantially increase the amount of product achieved (Tsai and Olson (1992) loc. cit.). Gelfiltration of contaminated DNA raw extracts has been used to physically separate DNA from inhibitors based on size differences (Tsai and Olson (1992) loc. cit.), (Jackson et. al. (1997) *Applied and Environmental Microbiology* 63, 4993-4995), (Miller (2001) *J Microbiol Methods* 44, 49-58). Charge differences between DNA and inhibitors were exploited in strategies using ion-exchange chromatography purification (Tebbe and Vahjen, (1993) loc. cit.); (Straub et. al. (1995) *Water Science and Technology* 31, 311-315); (Smalla et. al. (1993) *J Appl Bacteriol*, 74, 78-85). In a different approach substances showing selective affinity towards polyphenols like soluble polyvinylpyrrolidone (PVP, figure 2, relative molecular weight 10'000-360'000 Da), insoluble polyvinylpolypyrrolidone (PVPP, a crosslinked derivative of PVP) or CTAB (hexadecyltrimethylammonium bromide) have been used to absorb (Holben et. al. (1988) *Appl. Environ. Microbiol.* 54, 703-711) or precipitate inhibitors from solutions (Zhou et. al. (1996) loc. cit.). Berthelet and co-workers used a PVPP affinity-matrix to chromatograph contaminated DNA solutions on spin columns (Berthelet et. al. (1996) *FEMS Microbiol Lett* 138, 17-22). Using ultracentrifugally generated CsCl density gradients Holben and co-workers (Holben et. al. (1988) loc. cit.) purified DNA from inhibitors based on equilibrium densities. For the construction of high quality libraries of uniform and particularly large DNA insert sizes (in vectors like Cosmid, Fosmid, BAC) a high resolution size selection step is essential to provide the reaction with uniformly sized insert DNA, especially if like in the case for BACs the cloning process does not feature any inherent size selective steps. This makes gel electrophoresis particularly attractive for the purification of environmental DNA. Hereby charge-mass ratios and size differences can be exploited simultaneously to achieve kinetic resolution of DNA from inhibitors and simultaneously the DNA itself is spread out according to size. Although simple gel electrophoresis may suffice to produce clonable DNA from soils containing only small amounts of humic and fulvic acids (Rondon et. al. (2000) loc. cit.), the humic content of soils varies greatly and can reach up to 60-80 % of the total organic matter (Tsai and Rochelle (2001) *Environmental Molecular Microbiology*, Horizon Scientific Press, page 15-30 (Extraction of nucleic acids from environmental samples)). Mostly therefore additional purification steps are necessary and still failures to produce

clonable soil DNA are common (Entcheva et. al. (2001) Appl Environ Microbiol 67, 89-99).

A particular modification of this method was devised by Young and co-workers (Young et. al. (1993) loc. cit.). Here, PVP was added to the gel to selectively lower the charge-mass ratio of humic acids so to improve resolution from DNA. This technique combines an affinity-based selective purification step retarding inhibitors in an electric field with a DNA size resolution step that is indispensable in the preparation of insert DNA for efficient large fragment cloning in vectors like BAC.

Yet, prior to subsequent enzymatic modifications of environmental DNA as required for cloning (like PCR using e.g. *Taq* or *Pfu* DNA polymerases or fill-up reactions using e.g. *Klenow*- or T4-DNA-polymerase or ligations using e.g. T4-DNA-ligase or multi-step reactions like phage-packaging) absorbants like PVP must be removed as they themselves are inhibitory for further enzymatic processes.

The separation of agarose gel-purified DNA from the PVP absorbant is achieved in the prior art by employing affinity chromatography after melting the DNA containing agarose slice (GeneClean® in (Young et. al. (1993) loc. cit.)). Such a procedure, however, is not suitable to purify very large DNA molecules because shearing forces generated during the elution process will cause fragmentation. Additionally elution efficiency is inversely correlated with molecule size, so large molecules will be selectively lost.

Alternatively, large DNA molecules can be electroeluted from an agarose slice cut from a gel after electrophoresis. The DNA will be recovered in diluted form in a buffer like TAE and has to be concentrated before further manipulation. This routinely involves precipitation in 70% ethanol (Ausubel et. al. (Eds.) (1998) Current Protocols in Molecular Biology, John Wiley & Sons, 2.11-2.1.10). Yet, such an alcohol precipitation involves at least one further centrifugation step and, accordingly, adverse shearing forces. Consequently, in order to purify large DNA fragments from agarose gels for enzymatic manipulation and cloning, prior art procedures involve melting a gel slice containing DNA, performing in-gel enzymatic manipulations in a re-solidified gel (like end polishing, ligation), solubilizing the gel using a degrading enzyme (Gelase® Epicentre, USA) and transforming the DNA into hosts (*E. coli*).

These procedures are complex and may lead to fragmentation or even loss of nucleic acid molecules. Furthermore, these manipulations of the prior art can not be carried out in the presence of enzyme inhibiting substances like PVP.

Many of the above strategies in different combinations have been used as part of multi-step purification protocols to produce clonable "metagenome DNA". Yet, environmental DNA purification is not trivial and whereas purification of DNA for PCR purposes may be accomplished using commercial kits (FastPrep® *Bio101*, USA), the preparation of sufficient amounts of concentrated and inhibitor-free high-molecular weight DNA (20-300 kbp) for cloning into Cosmids or BACs is much more challenging and may be doomed with failure (Entcheva et. al. (2001) loc. cit).

Thus, the technical problem underlying the present invention was to provide means and method for cloning of genetic material isolated from primary samples. The solution to this technical problem is achieved by providing the embodiments characterized in the claims.

The current invention provides means to overcome the technical difficulties associated with the isolation and cloning of large fragment DNA from uncultivated environmental sources.

Accordingly, the present invention relates to a device for the isolation and/or purification of nucleic acid molecules comprising at least two layers, a first layer being adapted to bind and/or inactivate inhibitors of the activity of reagents or enzymes used in nucleic acid manipulation and a second layer being adapted to separate a plurality of nucleic acid molecules with respect to their size.

The term "device" as employed herein is an arrangement/construction comprising, inter alia, gels, or gel chambers or columns as defined herein. Preferably, the gels, gel chambers or columns form the device of the present invention. The nucleic acid molecules to be isolated and/or purified are isolated and/or purified by passing them through the at least two layers of the device as defined herein.

The term "inhibitors of the activity of reagents or enzymes used in nucleic acid manipulation" describes substances comprised in samples of soil, aquatic samples

or samples of symbiotic/parasitic consortia which inhibit the activity of reagents or enzymes used in nucleic acid manipulation. Examples of said substances are described herein above and comprise inorganic inhibitors, like, e.g. heavy metal ions, organic inhibitors, like polysaccharides and in particular humic and fulvic acids. Humic and fulvic acids are high molecular weight heterocyclic polyphenols mainly of plant origin with an affinity to polynucleotides and strongly protein denaturing properties. The chemical structure of humic and fulvic acids is shown in appended figures 1A and 1B. Moreover, the chemical properties of said groups of molecules is described in detail by Stevenson (Humus chemistry: genesis, composition, reactions (1994) Wiley New York) and Buffle (Les substances humiques et leurs interactions avec les ions minéraux (1977) Conference Proceedings de la Commission d'Hydrologie Appliquée de A.G.H.T.M.. l'Université d'Orsay, 3-10).

The term "reagents used in nucleic acid manipulation" as used in this context comprises substances like metal ions (e.g.  $Mg^{2+}$ ,  $Mn^{2+}$ ,  $Ca^{2+}$ ), (charged) inorganic and organic molecules required for enzymatic activity or for enzymatic co-factors, said co-factors themselves, or stabilizers. The term "enzymes used in nucleic acid manipulations" relates to enzymes like RNase(s), DNase(s), DNA-polymerase(s), ligase(s) or kinase(s) which are used for nucleic acid manipulation.

The term "nucleic acid manipulation" as used herein comprises standard methods known by the person skilled in the art. Said methods comprise DNA-engineering, such as cloning methods of nucleic acid molecules, the mutation of nucleotide sequences of nucleic acid molecules or amplification methods, like, e.g. PCR. Examples for said methods are described in the appended examples and in laboratory manuals, e.g. Sambrook et. al. (1989) *Molecular Cloning: A Laboratory Manual*, Second Edition, Cold Spring Harbor Laboratory Press, Cold Spring Harbor, New York; Ausubel et. al. (1998), loc. cit.. In particular, the term "nucleic acid manipulation" relates to the manipulation of DNA or RNA and corresponding cloning techniques.

The term "layer" defines in accordance with the present invention a physical matrix which is characterized by its ability to separate samples containing nucleic acid



molecules and also characterized by its ability to separate, if desired, different nucleic acid molecules by their physiological properties, like size or overall charge. The recited first layer is adapted to characterized by its ability to bind and/or inactivate inhibitors described herein above. The ability of the first layer to bind and/or inactivate inhibitors may be achieved by the addition of compounds with sufficiently high binding affinity to the above described inhibitors so to retard their mobility in aqueous solutions and reduce their effective free concentration so to relieve nucleic acids migrating through the device of the invention from comigrating inhibitors.

The recited second layer is characterized by its ability to separate a plurality of nucleic acid molecules with respect to their size. Accordingly, said physical matrix may be a form of a physical matrix, suited to separate nucleic acid molecules based on molecular sieving, e.g. comprising gels or polymers.

According to a preferred embodiment of the device of invention said first layer is arranged above the second layer.

The term "above" defines the position of the first layer relative to the second layer and relative to the direction in which the samples migrate through the layers of the device. Accordingly, the nucleic acid molecule to be isolated and/or purified with the device of the invention is first contacted with the physical matrix of the first layer and afterwards contacted with the matrix of the second layer. Thus, the present invention comprises devices in which the first layer is horizontally above the second layer as well as devices in which the first layer is vertically above the second layer. Such a device is illustratively exemplified in the appended examples as a gel comprising two phases/layers. Accordingly, as illustrated in the appended figures and examples and described herein, this device may be arranged in form of a gel. Therefore, the device of the invention may be a device, wherein said first layer is a first phase of a gel and said second layer is a second phase of said gel. In its broadest sense, the term "phase" of a gel indicates that this phase has a different overall chemical constitution than a further (second) phase of said gel. For example, the difference in the overall chemical construction may be due to the presence of chemicals/compounds that bind or inactivate the above mentioned inhibitors of the activity of reagents or enzymes used in nucleic acid manipulations.

Preferably, the device of the present invention comprises a gel, wherein said gel is an agarose-gel or a polyacrylamid-gel.

Methods for the preparation of said gels are known by the person skilled in the art and are described in the appended examples and in standard laboratory manuals, e.g. Sambrook et. al. (1989), loc. cit., Cold Spring Harbor, New York; Ausubel et. al. (1998), loc. cit..

Preferably, the device of the invention comprises in said first layer polyvinylpyrrolidone (PVP), or polyvinylpolypyrrolidone (PVPP) or combinations thereof. As demonstrated in the appended examples PVP and PVPP are immobile components of said first layer or are characterized at least by a lower mobility compared to the nucleic acid molecules and to bind or interact with the above characterized "inhibitors of the activity of reagents or enzymes used in nucleic acid manipulation".

Further examples of corresponding components of said first layer (i.e. inactivators of the inhibitors defined herein above) are functional molecules like CTAB, EDTA, EGTA, cyclodextrins, proteins, (poly)peptides, nucleic acids immobilized or tethered on appropriate matrices acting as catcher molecules or ion-exchanger. These functional molecules may act by complexing ions (like EDTA, EGTA,), precipitating polysaccharides (like CTAB), binding small hydrophobic molecules or uncharged small molecules, like cyclodextrins, binding specific surface structures (like proteins, (poly)peptides and aptamers). Said proteins may comprise, e.g., antibodies and (poly)peptides directed against inhibitors. Furthermore, lectins are envisaged as inactivators of the inhibitors defined herein. An example for a nucleic acids acting as catcher molecules in accordance with the invention are RNA-aptamers.

The term "(poly)peptide" as used herein summarizes a group of molecules which comprise the group of peptides, consisting of up to 30 amino acids, as well as the group of polypeptides, consisting of more than 30 amino acids.

Above described low molecular weight compounds comprise electrically charged compounds (e.g. CTAB, EDTA) and uncharged small molecules (e.g. cyclodextrin). Said compounds are soluble in aqueous solutions and, thus, mobile due to rapid diffusion and in particular mobile in an electrical field. In order to achieve a lower

mobility of said compounds compared to the nucleic acid molecules, said compounds may be coupled to the physical matrices (e.g. chemically coupled). Compounds which are soluble in aqueous solutions but not electrically charged (e.g. PVP) and compounds which are non-soluble in aqueous solutions (e.g. PVPP) show a lower mobility compared to the nucleic acid molecules and thus do not essentially require to be coupled to the physical matrices.

More preferably, the device the invention is a device, wherein said second layer is substantially free of PVP, PVPP, CTAB, EDTA, EGTA, cyclodextrins, proteins, (poly)peptides, nucleic acids or ion-exchanger.

The term "substantially free" is understood in accordance with the invention to define a layer which does not contain the recited compounds in an effective amount which is detectable by standard methods. Said standard methods are known in the art and comprise MS (mass spectrometry), FT-IR (fourier transform infrared spectrometry), NMR (nuclear magnetic resonance) or HPLC (high performance liquid chromatography). Accordingly, the second layer does, most preferably, not contain any of the recited compounds as an essential/effective element.

According to a more preferred embodiment the device of the invention is electrically biased to enhance flow of (a) sample(s) through the layers.

As known by the person skilled in the art nucleic acid molecules are negatively charged due to the ribose-phosphate framework. Accordingly, nucleic acid molecules migrate in an electrical field from the cathode (-) to the anode (+).

Examples for devices which are electrically biased to enhance flow of (a) sample(s) are devices for gel electrophoresis. Devices with continuous as well as devices with discontinuous electrical fields are particularly comprised by the present invention. Accordingly, the electrical field can be discontinuous due to a gradient of a salt (buffer salt), a pulsed field comprising e.g. different angles. Again, an electrically biased device in accordance with the present invention is shown in the appended figures and illustrated in the appended examples.

The invention relates to a device, wherein said first layer preferably comprises sample loading means.

The term "sample loading means" defines in accordance with the invention means for placing the sample comprising nucleic acid molecules in the device of the invention. Examples for said means are sample slots in a gel, the surface of the matrix of a column or a valve for injecting a sample onto a column.

The isolation of nucleic acid molecules by using a device of the invention in the format of a horizontal agarose-gel is described in the appended example 1 and in the format of a single column is described in the appended example 4. Examples for the corresponding devices are shown in the appended figures 3 and 4.

In a more preferred embodiment of the invention said loading means are provided in an array in an upper portion of the first layer, defining an array of columns, each being capable of isolating nucleic acid molecules.

Hence, a device according to the preferred embodiment of the invention comprises more than one means for placing the sample. Thus, it is possible to isolate nucleic acid molecules from different samples in parallel. An example for loading means provided in an array is a gel comprising different loading slots and, therefore, different lanes (lines in which the samples are separated). An alternative example for said array are columns which are arranged in groups, e.g. in a frame. Said frame may comprise low numbers of columns (two to twelve) for the isolation of nucleic acid molecules from low numbers of samples. Also in accordance with the invention are frames comprising medium numbers of columns (twelve to 96) as well as frames with high numbers of columns (more than 96) which are suitable for high throughput screens (HTS).

According to an alternative embodiment of the invention said first layer of the device is arranged below the second layer.

The term "below" defines the position of the second layer relative to the first layer and relative to the direction in which the samples migrate through the layers of the device. The present invention comprises devices for the isolation and/or purification of nucleic acid molecules in which the nucleic acid molecules are first contacted with the physical matrix of the second layer to separate the molecules with respect to their size and subsequently contacted with the physical matrix of the first layer to bind and/or inactivate the above defined inhibitors.

Accordingly, the present invention comprises devices in which the second layer is horizontally above the first layer as well as devices in which the second layer is vertically above the first layer. Most preferred is in this context a device in form of a column, wherein said column comprises said first and said second layer.

An alternatively preferred embodiment of the invention relates to a device, wherein said second layer is a first phase of a column and said first layer is a second phase of said column. As pointed out herein above, the first layer comprises functional molecules capable of inactivating or binding inhibitors of the activity of reagents or enzymes used in nucleic acid manipulation.

It is also envisaged, in accordance with the invention, that an enhancement of the flow of (a) sample(s) through the layers of the device is, for example, achieved by gravity or by the pressure of the flow of a diluent. Examples for appropriate diluents are buffer solutions. The appended Examples show a device in form of a column as described herein.

Preferably, the device of the invention comprises in said first layer (in a column preferably the second phase) a matrix comprising PVP or PVPP or combinations thereof. As described herein above PVP and PVPP are characterized at least by a lower mobility compared to the nucleic acid molecules and to bind or interact with the above characterized "inhibitors of the activity of reagents or enzymes used in nucleic acid manipulation".

Further examples of corresponding components of said first layer (i.e. compounds capable of inactivating the inhibitors of nucleic acid manipulations) are EDTA, EGTA, CTAB, cyclodextrins, proteins, (poly)peptides, nucleic acids acting as catcher molecules or ion-exchangers. Said proteins may comprise, e.g., antibodies directed against inhibitors. Also lectins are envisaged as inactivators. An example for said nucleic acids acting as catcher molecules in accordance with the invention are RNA-aptamers.

More preferably the device the invention is a device, wherein said second layer (in a column preferably the first phase) is a matrix which is substantially free of PVP,

PVPP, CTAB, EDTA, EGTA, cyclodextrins, proteins, (poly)peptides, nucleic acids or ion-exchangers.

The term "substantially free" is defined herein above.

According to a further preferred embodiment of the invention, said matrix of said first and/or second layer is agarose, sepharose<sup>TM</sup>, sephadex<sup>TM</sup>, sephacryl<sup>TM</sup>, BioGel<sup>TM</sup>, superose<sup>TM</sup> or acrylamide.

Variations of the samples comprising the nucleic acid molecules to be isolated and/or purified with the device of the invention may require specific materials of the matrix of the first and/or the second layer. Said variations may concern the quality of the comprised nucleic acid molecules as well as the quality of the sample itself with respect to characteristic substances which may be contained. Accordingly, the matrix of the first and/or the second layer may be a specific matrix for gel-filtration which allows a molecular sieving of the nucleic acid molecules. Suitable media for matrices are characterized by a specific size of the pores which is known by the person skilled in the art and described in the instructions provided by the manufacturer of commercially available media. Said matrices comprise media for gel-filtration. Examples of said media comprise sepharose<sup>TM</sup> (e.g. sepharose2B, sepharose4B, sepharose6B), sephadex<sup>TM</sup> (e.g. sephadex G200, sephadex G150) sephacryl<sup>TM</sup> and superose<sup>TM</sup> as well as bio-gel<sup>TM</sup> P100 and bio-gel<sup>TM</sup> P200. Moreover, said media comprise agarose, polymers as e.g. dextrans and acylamid based-resins. In particular, media are preferred which are suitable for two-phases-columns. Further suitable materials are known to the person skilled in the art.

In a further preferred embodiment of the device of the invention said nucleic acid molecule is DNA or RNA. Most preferably, said DNA is genomic DNA (gDNA).

A particularly preferred embodiment of the invention is a device, wherein said nucleic acid molecule is derived from (micro)organisms of soil, sediments, water, for example sea water, or symbiotic/parasitic consortia.

The term "soil" defines in accordance with the invention the complex product of geological and biological processes acting on inorganic minerals and biomass deposited on the earth surface. It contains the majority of microbial biodiversity on

earth (Whitman et. al. (1998) Proc Natl Acad Sci USA, 95(12, 6578-83) acting to recycle and biomineralize organic matter and serves as a substratum to anchor and nourish higher plants.

In the appended examples the preparation of nucleic acid molecules isolated from soil is exemplarily described.

Examples for symbiotic/parasitic consortia in accordance with the present invention are consortia isolated e.g. from animal tissues or organs, e.g. from gut, stomach, intestines, like appendix or insect-, bird- and mammalian-intestinal tracts or -guts, comprising ruminant-gut. Also envisaged are animal tissues or organs from annelid(s), mollusc(s), sponge(s), cnidaria, arthropod(s), amphibian(s), fish or reptile(s). However, it is also envisaged that nucleic acids from parasitic consortia from human tissue, organs, sputum, faeces, sperm, blood, urine or other body fluids are isolated.

More preferably said (micro)organisms from which said nucleic acid molecules are derived from are (micro)organisms of aquatic plankton, animal tissues and organs as described herein above, microbial mats, clusters, sludge flocs, or biofilms.

(Micro)organisms of the "aquatic plankton" comprises bacterial plankton, archaeal plankton, viruses, phytoplankton as well as zooplankton. Said (micro)organisms are known as small organisms.

Biofilms are microbial assemblages on a surface in "aqueous environments" in which microbes are embedded in a hydrated polymeric matrix. This matrix acts like a kind of glue, holding the microbes together, attaching them to the surface and protecting them from detrimental external influences. They may contain several taxonomically distinct species (e.g. bacteria, fungi, algae, and protozoa) and may form on surfaces of diverse composition like e.g. metals, glass, plastics, tissue, mineral and soil particles.

Microbial mats and clusters are microbial assemblages/aggregates similar to biofilms in composition however not necessarily as firmly attached to solid surfaces.

According to a further preferred embodiment of the invention said (micro)organisms are isolated and/or purified as consortia of coexisting species. Preferably it is envisaged that said (micro)organisms are isolated and/or purified as consortia of

coexisting species without previous separation/purification of single microorganismic species.

Preferably said consortia of coexisting species comprise at least one organism that cannot proliferate indefinitely in an artificial setting (e.g. a synthetic or semisynthetic culture medium) in the absence of other species and/or in the absence of the substratum it is isolated from, and wherein said substratum contains the above defined inhibitors of the activity of reagents or enzymes used in nucleic acid manipulation.

To obtain the DNA of said microorganisms they may either be bulk-separated mechano-chemically from most of the surrounding matrix they are attached to or embedded or suspended in (like water, soil, sediments, organic debris of plant or animal origin) before lysis (indirect lysis) or they may be lysed and extracted directly i.e. still in the context of/attached to their surrounding physical matrix (soil, biofilm, floc, cluster) that contains a plurality of potential inhibitors of molecular, in particular enzymatic manipulation. Bulk-separation of microorganisms may be accomplished e.g. through mechanical agitation, ion-exchange resin mediated desorption (optionally facilitated through substances added to the extraction buffers used like detergents, salts) followed by an optional concentration step like filtration (for suspended plankton) or suspending and differential gravitational sedimentation in a liquid. In both instances total nucleic acids of mixed origin are isolated irrespective of taxonomic status, abundance and cultivability of the respective taxonomically mixed species.

Classic isolation of nucleic acid molecules from single microbial species or groups of microbial species comprise cultivation of said organisms by massive dilution in selectively growth supporting media. Thereby above mentioned inhibitors of molecular manipulation of nucleic acids are diluted to facilitate subsequent isolation and cloning. Yet at the cost of a massive reduction in species representation as very few species can be supported by standard cultivation techniques (see herein above). In contrast the present invention provides means for the isolation of nucleic acid molecules derived from taxonomically mixed (micro)organisms without cultivation in synthetic media. Since such steps of cultivation under laboratory conditions results in depletion or at least in significant dilution of organisms which do not grow under said



conditions the present invention surprisingly provides means for the isolation and/or purification of nucleic acid molecules derived from said organisms.

In a preferred embodiment of the invention said nucleic acid molecules represent a fraction of the metagenome of a given habitat.

As known in the art the term "metagenome" defines the totality of all genomes of organisms of a given habitat and is furthermore defined in the art; see inter alia Handelsman et. al. (1998) loc. cit.. In particular, the term "metagenome" relates to genomic nucleic acids, preferably DNA, derived from unknown or uncultivable microorganisms, i.e. organisms that cannot be isolated by standard methods and made actively replicating in standard artificial media for indefinite periods of time.

Accordingly the term "a fraction of the metagenome of a given habitat" defines in accordance with the invention nucleic acid molecules and in particular large nucleic acid molecules (>200bp) derived from the total pool of heterogenous microbial genomes present in a given habitat. This is irrespective of phylogenetic affiliation or molecular or physiological traits. Particularly the representation of any particular microbial genome in the extracted portion of the metagenome is not influenced by or dependent on the cultivatability of this organism. Therefore nucleic acids of uncultivated and in a preferred form particularly of uncultivable (micro)organisms are substantially represented in the extracted fraction of the metagenome.

An alternative embodiment of the invention relates to a method for the isolation of a nucleic acid molecule comprising applying a sample to the device as defined herein above. Said sample (for example derived from soil, sediments, water or symbiotic/parasitic consortia) is loaded onto a device comprising the above-identified at least two layers (for example a gel or a column comprising said two layers) and the nucleic acid molecule is purified and/or isolated by passing them through said two layers. In accordance with the invention, the first layer as defined herein, i.e. the layer comprising the inactivators of the inhibitors of reagents and enzymes of nucleic acid manipulations,4 retains said inhibitors, wherein said second layer provides for, e.g., an isolation and separation step for isolating/separating the nucleic acid molecules in accordance with their physical properties, like size or overall charge. It is envisaged and documented in the appended examples, that the nucleic acid

molecules to be isolated pass completely through the two layers of the device of the invention (e.g. a column comprising the two layers) or that the nucleic acid molecules pass only partially through the second layer. For example, it is envisaged that said nucleic acid molecules pass completely through said first layer comprising the inactivators of inhibitors and only partially through said second layer which is substantially free of said inactivators. This method is, inter alia, employed in the device of the invention which is in form of a gel. The nucleic acid molecules to be isolated and/or purified may be isolated or purified from the second layer by conventional means, for example by electroelution.

The method of the invention is also exemplified in the appended examples 1, 4 and 5.

Said method may optionally comprise one or more additional steps of subsequent purification of the obtained nucleic acid molecule(s).

According to a preferred embodiment of the method of the invention a fractions of the metagenome is isolated from a given habitat.

The invention relates in an alternative embodiment to a method for the generation of at least one gene library, comprising the steps of

- (a) isolating and/or purifying nucleic acid molecules from a sample using a device as defined herein above and optionally amplifying said nucleic acid molecules;
- (b) cloning the isolated and/or purified and optionally amplified nucleic acid molecules into appropriate vectors; and
- (c) transforming suitable hosts with said suitable vectors.

Methods for the amplification of an isolated and/or purified nucleic acid molecule are known in the art and comprise e.g. polymerase chain reaction (PCR).

Methods for the cloning of nucleic acid molecules into appropriate vectors and transforming suitable hosts with said suitable vectors represent standard methods of molecular biology and are described in the appended examples (in particular, example 2) and in laboratory manuals, e.g. Sambrook et. al. (1989) loc. cit.; Ausubel et. al. (1998), loc. cit.; Mülhardt (2002) Der Experimentator: Molekularbiologie/Genomics; Gustav Fischer. Suitable vectors are described herein below.

The above described method may, optionally, additionally comprise one or more of the following steps prior to the cloning of nucleic acid molecules into appropriate vectors according to step (b):

- (i) modification of DNA ends of the isolated and/or purified nucleic acid molecules, e.g. to remove or fill-up random 3'- or 5'-overhangs (polishing/fill-up/blunting), by DNA polymerase treatment (e.g. with Klenow enzyme, T4-polymerase) or treatment with exonucleases (like mung bean nuclease) or introducing defined 3'overhangs (like single nucleotide overhangs, in particular adenosin overhangs) using, e.g. DNA polymerase (like Taq), for subsequent cloning into appropriate vectors (e.g. T-overhang vectors like pGEM T-easy from Promega, USA) or introducing 3'- or 5'-overhangs using restriction endonucleases;
- (ii) phosphorylation of the isolated and/or purified nucleic acid molecules (e.g. by PNK); and/or
- (iii) ligation of the isolated and/or purified nucleic acid molecules to other nucleic acid molecules by treatment with an enzymatic ligase (e.g. by T4-ligase) or topoisomerase.

Also envisaged is a "sizing" step, as also illustrated briefly in (i) herein above, wherein said step comprises sizing the obtained nucleic acid molecules by treatment of the isolated and/or purified nucleic acid molecules with an enzymatic endonuclease (e.g. by restriction endonucleases or DNase I) and/or mechanical shearing (e.g. by ultrasonication or passing nucleic acids with high pressure through narrow tubes or valves similar to the "nebulizer" from Invitrogen, USA).

Suitable vectors comprise plasmid vectors (e.g. pUC18 and derivates thereof, pBluescript etc.), cosmid vectors (e.g. Expand, SuperCos), fosmid vectors (e.g. EpiFos 5), phage vectors (e.g. lambda ZAP), BAC vectors (e.g. pBeloBAC) and YAC vectors.

Preferably, said suitable hosts are selected from the group consisting of *E. coli*, *Pseudomonas* sp., *Bacillus* sp, *Streptomyces* sp, other actinomycetes, myxobacteria yeasts and filamentous fungi.

Transformation of said suitable host may be assayed by standard methods of molecular biology; see Sambrook et. al. (1989) loc. cit.; Ausubel et. al. (1998), loc.

cit.; Mülhardt (2002) loc. cit.. Corresponding assays for a successful transformation may be based on sequence similarity and performed on plated colonies (filter hybridization) by probe hybridization or by PCR analysis of colonies or extracted DNA of single or multiple colonies. Methods for the preparation of suitable oligonucleotides are known in the art. Alternatively the inserts may be sequenced and targets identified via homology search in appropriated sets of deposited sequence data (e.g. GenBank).

Activity based assays may be performed by screening for substrate conversion/degradation inside or outside the host (scoreable by substrate clearing zones around colonies; color development; molecular product profiling by high performance liquid chromatography (HPLC), mass spectrography (MS) or gas chromatography (GC); complementation of growth-deficient mutants), by screening for growth inhibition/stimulation of indicator organisms (in overlays).

Accordingly, and in a further embodiment, the present invention relates to a gene library obtained by the method disclosed herein and by employing the device described in this invention.

An alternative embodiment of the invention relates to a gene library generated from metagenome nucleic acid molecules, preferably from DNA, from non-planctonic (micro)organisms comprising average insert sizes of at least 50 kB, at least 55 kB, at least 60 kB, at least 70 kB, at least 80 kB, at least 90 kB or at least 100 kB.

As defined herein above planctonic (micro)organisms of the "aquatic plancton" comprises bacterial- and archaeal plancton, viruses, phytoplankton as well as zooplankton. Said (micro)organisms are known as small organisms living in aquatic habitats. Accordingly the term "non-planctonic (micro)organisms" defines (micro)organisms in accordance with the invention which are not comprised by the term "planctonic (micro)organisms". This group of (micro)organisms comprise (micro)organisms of soil, microbial mats, clusters sludge flocs, biofilms and symbiotic/parasitic consortia.

Similarly, the invention relates to a gene library generated from metagenome nucleic acid molecules, preferably DNA, from planctonic (micro)organisms comprising

average insert sizes of at least 85 kB, at least 90 kB, at least 95 kB, at least 100 kB, at least 120 kB, at least 140 kB, at least 160 kB or at least 200 kB.

In contrast to the group of "non-planctonic (micro)organisms" defined herein above the term "planctonic (micro)organisms" defines (micro)organisms of the "aquatic plancton" comprises bacterial- and archaeal plancton, viruses, phytoplankton and zooplankton as described herein above.

The average insert size of a (gene) library is, inter alia, determined by (a) isolating the cloned recombinant DNA of at least 0.1% of the clones of the respective library (however no less than a minimum of 20 clones) by methods known in the art, (b) digesting the isolated cloned recombinant DNA molecules with restriction enzymes (6-base or 8-base cutters, e.g., BamHI or EcoRI or Not I used singular or in combination) so to preferentially digest the vector backbone away from insert DNA (e.g. Not I used for pEpiFos5), (c) separating the resulting DNA fragments obtained from each clone individually by agarose gel electrophoresis (continuous or pulsed-field) as known in the art and (d) adding together the sizes of all non-vector fragments of all analyzed clones of a library and dividing the resulting number by the number of clones analyzed in order to obtain a figure for the average insert size.

The present invention relates further to a nucleic acid molecule comprising a DNA as depicted in SEQ ID NO: 1 or comprising a DNA as deposited under EMBL accession number AJ496176.

Said nucleic acid molecule of the invention has been isolated and obtained by employing the device of the invention. It represents a part of the genome of the newly identified Crenarchaeote as isolated with methods described herein and by techniques implied in the device of the present invention. Taxonomically the crenarchaeota represent a prokaryotic phylum as part of the archeal kingdom. The majority of its representatives are known to be hyperthermophiles yet increasingly they are found in mesophilic habitats as well; see Burggraf et al. (1997) Int J Syst Bacteriol., 47, 657-660; Preston et al. (1996) Proc Natl Acad Sci USA., 93, 6241-46. In the context of the present invention, the term "genome" defines not only sequences which are open reading frames (ORFs) encoding proteins, polypeptides

or peptides, but also refers to non-coding sequences. Accordingly, the term "nucleic acid molecule" comprises coding and, wherever applicable, non-coding sequences. The nucleic acid molecule of the invention furthermore comprises nucleic acid sequences which are degenerative to the above nucleic acid sequences. In accordance with the present invention, the term "nucleic acid molecule" comprises also any feasible derivative of a nucleic acid to which a nucleic acid probe may hybridize. Said nucleic acid probe itself may be a derivative of a nucleic acid molecule capable of hybridizing to said nucleic acid molecule or said derivative thereof. The term "nucleic acid molecule" further comprises peptide nucleic acids (PNAs) containing DNA analogs with amide backbone linkages (Nielsen, P., Science 254 (1991), 1497-1500). The term "ORF" ("open reading frame") which encodes a polypeptide, in connection with the present invention, is defined either by (a) the specific nucleotide sequences encoding the polypeptides specified above in (aa) or in (ab) or (b) by nucleic acid sequences hybridizing under stringent conditions to the complementary strand of the nucleotide sequences of (a) and encoding a polypeptide deviating from the polypeptide of (a) by one or more amino acid substitutions, deletions, duplications, insertions, recombinations, additions or inversions.

Furthermore the present invention relates in one embodiment to a nucleic acid molecule representing part of the genome of a non-thermophilic crenarchaeote, whereby said nucleic acid molecule has at least one of the following features:

- (a) it contains at least one ORF which encodes a polypeptide having the amino acid sequence SEQ ID NO: 3, SEQ ID NO: 5, SEQ ID NO: 7, SEQ ID NO: 9, SEQ ID NO: 11, SEQ ID NO: 13, SEQ ID NO: 15, SEQ ID NO: 17, SEQ ID NO: 19, SEQ ID NO: 21, SEQ ID NO: 23, SEQ ID NO: 25, SEQ ID NO: 27, SEQ ID NO: 29, SEQ ID NO: 31, SEQ ID NO: 33, SEQ ID NO: 35;
- (b) comprises the DNA sequence of SEQ ID NO: 1, SEQ ID NO: 2, SEQ ID NO: 4, SEQ ID NO: 6, SEQ ID NO: 8, SEQ ID NO: 10, SEQ ID NO: 12, SEQ ID NO: 14, SEQ ID NO: 16, SEQ ID NO: 18, SEQ ID NO: 20, SEQ ID NO: 22, SEQ ID NO: 24, SEQ ID NO: 26, SEQ ID NO: 28, SEQ ID NO: 30, SEQ ID NO: 32, SEQ ID NO: 34;

- (c) it comprises portion of at least 20 nucleotides, preferably 100 nucleotides, more preferably at least 500 nucleotides which hybridize under stringent conditions to the complementary strand of SEQ ID NO: 2, SEQ ID NO: 4, SEQ ID NO: 6, SEQ ID NO: 8, SEQ ID NO: 10, SEQ ID NO: 12, SEQ ID NO: 14, SEQ ID NO: 16, SEQ ID NO: 18, SEQ ID NO: 20, SEQ ID NO: 22, SEQ ID NO: 24, SEQ ID NO: 26, SEQ ID NO: 28, SEQ ID NO: 30, SEQ ID NO: 32, SEQ ID NO: 34;
- (d) it is degenerate as a result of the genetic code with respect to the nucleic acid molecule of (c); or
- (e) it is at least 50% identical with the nucleic acid molecule of SEQ ID NO: 2, SEQ ID NO: 20 or SEQ ID NO: 30, 45% identical with the nucleic acid molecule of SEQ ID NO: 8 or SEQ ID NO: 26, 35% identical with the nucleic acid molecule of SEQ ID NO: 16, SEQ ID NO: 22 or SEQ ID NO: 24 or 30% identical with the nucleic acid molecule of SEQ ID NO: 4, SEQ ID NO: 14, SEQ ID NO: 18 or SEQ ID NO: 28;

The following list relates to the SEQ ID NOS as defined herein and shows (partial) identification of ORFs as defined herein:

Name	Identity
SEQ ID NO:1	complete DNA sequence 29i4, 33925 nt, of a newly identified Crenarchaeote
SEQ ID NO:2	ORF001 DNA sequence, 2367 nt, Fam. B DNA Polymerase (truncated ORF)
SEQ ID NO:3	ORF001 Protein sequence, 789 aa, Fam. B DNA Polymerase (truncated protein)
SEQ ID NO:4	ORF002 DNA sequence, 882 nt, $\alpha/\beta$ hydrolase
SEQ ID NO:5	ORF002 Protein sequence, 293 aa, $\alpha/\beta$ hydrolase
SEQ ID NO:6	ORF003 DNA sequence, 318 nt
SEQ ID NO:7	ORF003 Protein sequence, 105 aa
SEQ ID NO:8	ORF004 DNA sequence, 1086 nt, Polyhydroxyalkanoate Synthase
SEQ ID NO:9	ORF004 Protein sequence, 361 aa, Polyhydroxyalkanoate

	Synthase
SEQ ID NO:10	ORF005 DNA sequence, 582 nt
SEQ ID NO:11	ORF005 Protein sequence, 193 aa
SEQ ID NO:12	ORF006 DNA sequence, 438 nt
SEQ ID NO:13	ORF006 Protein sequence, 145 aa
SEQ ID NO:14	ORF007 DNA sequence, 915 nt, Glycosyl Transferase group 1
SEQ ID NO:15	ORF007 Protein sequence, 304 aa, Glycosyl Transferase group 1
SEQ ID NO:16	ORF008 DNA sequence, 1692 nt, Asparagine Synthase
SEQ ID NO:17	ORF008 Protein sequence, 563 aa, Asparagine Synthase
SEQ ID NO:18	ORF009 DNA sequence, 666 nt, Phosphoserin Phosphatase
SEQ ID NO:19	ORF009 Protein sequence, 221 aa, Phosphoserin Phosphatase
SEQ ID NO:20	ORF010 DNA sequence, 1212 nt
SEQ ID NO:21	ORF010 Protein sequence, 403 aa
SEQ ID NO:22	ORF011 DNA sequence, 1164 nt, Transmembrane protein
SEQ ID NO:23	ORF011 Protein sequence, 387 aa, Transmembrane protein
SEQ ID NO:24	ORF012 DNA sequence, 882 nt, Fix A Electron Transfer Flavoprotein
SEQ ID NO:25	ORF012 Protein sequence, 293 aa, Fix A Electron Transfer Flavoprotein
SEQ ID NO:26	ORF013 DNA sequence, 1284 nt, Fix B Electron Transfer Flavoprotein
SEQ ID NO:27	ORF013 Protein sequence, 427 aa, Fix B Electron Transfer Flavoprotein
SEQ ID NO:28	ORF014 DNA sequence, 1878 nt, Fix CX Fusion Electron Transfer Flavoprotein
SEQ ID NO:29	ORF014 Protein sequence, 625 aa, Fix CX Fusion Electron Transfer Flavoprotein
SEQ ID NO:30	ORF015 DNA sequence, 2238 nt, Sensory Transduction Histidin Kinase
SEQ ID NO:31	ORF015 Protein sequence, 745 aa, Sensory Transduction Histidin Kinase
SEQ ID NO:32	ORF016 DNA sequence, 519 nt



SEQ ID NO:33	ORF016 Protein sequence, 172 aa
SEQ ID NO:34	ORF017 DNA sequence, 1008 nt, (truncated ORF)
SEQ ID NO:35	ORF017 Protein sequence, 335 aa, (truncated protein)

A potential field of application of "ORF004" as defined herein comprise the generation or modification of biogenic polymers/polyesters (polyhydroxyalkanoate synthase, Zinn, (2001) Adv Drug Deliv Rev 53(1):5-21 and Fidler, (1992), FEMS Microbiol Rev 9(2-4):231-5; Snell, (2002) Metab. Eng. 4(1):29-40)

Furthermore, the ORF008, a potential asparagine synthetase may be used in the context of amino acid synthesis (EC 6.3.5.4) and/or for the generation of transgenic organisms, like bacteria, plants with altered capacities to generate amino acids.

In addition, ORFs 12, 13 or 14 may play a role in redox processes involved in nitrogen fixation and may be useful in generating transgenic organisms like bacteria, plants with altered capacities to assimilate nitrogen.

The term "hybridizing" as used herein refers to a pairing of polynucleotides to a complementary strand of polynucleotide which thereby form a hybrid. Said complementary strand polynucleotides are, e.g. the polynucleotides of the invention or parts thereof. Therefore, said polynucleotides may be useful as probes in Northern or Southern Blot analysis of RNA or DNA preparations, respectively, or can be used as oligonucleotide primers in PCR analysis dependent on their respective size. Preferably, said hybridizing polynucleotides comprise at least 10, more preferably at least 15 nucleotides in length while a hybridizing polynucleotide of the present invention to be used as a probe preferably comprises at least 100, more preferably at least 200, or most preferably at least 500 nucleotides in length.

It is well known in the art how to perform hybridization experiments with nucleic acid molecules, i.e. the person skilled in the art knows what hybridization conditions s/he has to use in accordance with the present invention. Such hybridization conditions are referred to in standard text books such as Sambrook et. al. (1989) loc. cit. or Higgins, S.J., Hames, D. "RNA Processing: A practical approach", Oxford University

Press (1994), Vol. 1 and 2.

"Stringent hybridization conditions" (also referred to highly stringent conditions as contrasted to conditions of low stringency) refers to conditions which comprise, e.g. an overnight incubation at 42°C in a solution comprising 50% formamide, 5x SSC (750 mM NaCl, 75 mM sodium citrate), 50 mM sodium phosphate (pH 7.6), 5x Denhardt's solution, 10% dextran sulfate, and 20 µg/ml denatured, sheared salmon sperm DNA, followed by washing the filters in 0.1 x SSC at about 65°C. Said conditions for hybridization are also known by a person skilled in the art as "high stringent conditions for hybridization". Also contemplated are nucleic acid molecules that hybridize to the polynucleotides of the invention at lower stringency hybridization conditions ("low stringent conditions for hybridization"). Changes in the stringency of hybridization and signal detection are primarily accomplished through the manipulation of formamide concentration (lower percentages of formamide result in lowered stringency); salt conditions, or temperature. For example, lower stringency conditions include an overnight incubation at 37°C in a solution comprising 6X SSPE (20X SSPE = 3M NaCl; 0.2M NaH<sub>2</sub>PO<sub>4</sub>; 0.02M EDTA, pH 7.4), 0.5% SDS, 30% formamide, 100 µg/ml salmon sperm blocking DNA; followed by washes at 50°C with 1 X SSPE, 0.1% SDS. In addition, to achieve even lower stringency, washes performed following stringent hybridization can be done at higher salt concentrations (e.g. 5X SSC). Note that variations in the above conditions may be accomplished through the inclusion and/or substitution of alternate blocking reagents used to suppress background in hybridization experiments. Typical blocking reagents include Denhardt's reagent, BLOTTO, heparin, denatured salmon sperm DNA, and commercially available proprietary formulations. The inclusion of specific blocking reagents may require modification of the hybridization conditions described above, due to problems with compatibility.

Preferred in accordance with the present inventions are polynucleotides which are capable of hybridizing to the polynucleotides of the invention or parts thereof, under stringent hybridization conditions, i.e. which do not cross hybridize to unrelated polynucleotides.

The nucleic acid molecules that are homologous to the above-described molecules and that represent derivatives of these molecules usually are variations of these

molecules that represent modifications having the same biological function. They can be naturally occurring variations, for example sequences from other organisms, or mutations that can either occur naturally or that have been introduced by specific mutagenesis. Furthermore, the variations can be synthetically produced sequences. The allelic variants can be either naturally occurring variants or synthetically produced variants or variants produced by recombinant DNA processes.

Generally, by means of conventional molecular biological processes it is possible (see, e.g., Sambrook et. al. (1989) loc. cit.) to introduce different mutations into the nucleic acid molecules of the invention. One possibility is the production of deletion mutants in which nucleic acid molecules are produced by continuous deletions from the 5'- or 3'-terminus of the coding DNA sequence and that lead to the synthesis of proteins that are shortened accordingly. Another possibility is the introduction of single-point mutation at positions where a modification of the amino acid sequence influences, e.g., the enzyme activity or the regulation of the enzyme. By this method muteins can be produced, for example, that possess a modified  $K_m$ -value or that are no longer subject to the regulation mechanisms that normally exist in the cell, e.g. with regard to allosteric regulation or covalent modification. Such muteins may be identified, e.g. by methods of the present invention, to be valuable as therapeutically useful modulators (inhibitors/antagonists or enhancer/agonists) of the activity of the proteins of the present invention.

Nucleic acid molecules that hybridize to polynucleotides of the invention can be isolated, e.g., from genomic or cDNA libraries. In order to identify and isolate such nucleic acid molecules the polynucleotides of the invention or parts of these polynucleotides or the reverse complements of these polynucleotides can be used, for example by means of hybridization according to conventional methods (see, e.g., Sambrook (1989), loc. cit.). As a hybridization probe nucleic acid molecules can be used, for example, that have exactly or basically the nucleotide sequence of a part of the sequence shown in SEQ ID No: 1 or sequences complementary thereto. The fragments used as hybridization probe can be synthetic fragments that were produced by means of conventional synthesis methods and the sequence of which basically corresponds to the sequence of a nucleic acid molecule of the invention. Preferably, the nucleic acid molecule of the invention is DNA or RNA.

An alternative embodiment of the invention relates to a vector comprising an above defined nucleotide acid molecule.

The vector of the present invention may be, e.g., a plasmid, phagemid, cosmid, fosmid, BAC, virus, bacteriophage or another vector used e.g. conventionally in genetic engineering, and may comprise further genes such as marker genes which allow for the selection of said vector in a suitable hosts and under suitable conditions.

Furthermore, the vector of the present invention may, in addition to the nucleic acid molecule of the invention, comprise expression control elements, allowing proper expression of the coding regions in suitable hosts. Such control elements are known to the artisan and may include a promoter, a splice cassette, translation initiation codon, translation and insertion site for introducing an insert into the vector. Preferably, the nucleic acid molecule of the invention is operably linked to said expression control sequences allowing expression in eukaryotic or prokaryotic cells.

Many suitable vectors are known to those skilled in molecular biology, the choice of which would depend on the function desired and include plasmids, phagemid, cosmids, fosmid, BAC, virus, bacteriophages and other vectors used conventionally in genetic engineering. Methods which are well known to those skilled in the art can be used to construct various plasmids and vectors; see, for example, the techniques described in Sambrook (1989) loc. cit. and Ausubel (1998) loc. cit.. Alternatively, the nucleic acid molecule and vectors of the invention can be reconstituted into liposomes for delivery to target cells. Thus, according to the invention relevant sequences can be transferred into expression vectors where expression of a particular (poly)peptide/protein is required. Typical cloning vectors include pBscpt sk, pGEM, pUC9, pBR322 and pGBT9. Typical expression vectors include pTRE, pCAL-n-EK, pESP-1, pOP13CAT. Typical prokaryotic cloning and expression vectors include: plasmid vectors like the pUC series (e.g. pUC18, pUC19), pGEM series (e.g. pGEM 7zf+, Promega, USA), pET series (e.g. pET22B, Novagen, USA), pBBC 1MCS series, pNOF (GL Biotech Germany), pCR-TOPO series and pCR Blunt (Invitrogen, USA), pBluescript series, pCAL series and pBC series (Stratagene, USA); Fosmid vectors like pEpifos5 and pCC1 (Epicentre, USA); Cosmid vectors like the Expand series (Expand I, II, III) (Roche, Germany), SuperCos (Stratagene, USA), pOJ436; BAC vectors like pBeloBAC, pCC1BAC (Epicentre, USA). However, the

present invention also envisages the expression of nucleic acid molecules as disclosed herein in eukaryotic vectors. Preferably, said nucleic acid molecules are linked to "control sequences". Said linking may be direct or indirect and refers, preferably to an operable linkage.

The term "control sequence" refers to regulatory DNA sequences which are necessary to effect the expression of coding sequences to which they are ligated. The nature of such control sequences differs depending upon the host organism. In prokaryotes, control sequences generally include promoter, ribosomal binding site, and terminators. In eukaryotes generally control sequences include promoters, terminators and, in some instances, enhancers, transactivators or transcription factors. The term "control sequence" is intended to include, at a minimum, all components the presence of which are necessary for expression, and may also include additional advantageous components.

The term "operably linked" refers to a juxtaposition wherein the components so described are in a relationship permitting them to function in their intended manner. A control sequence "operably linked" to a coding sequence is ligated in such a way that expression of the coding sequence is achieved under conditions compatible with the control sequences. In case the control sequence is a promoter, it is obvious for a skilled person that double-stranded nucleic acid is preferably used.

Thus, the vector of the invention is preferably an expression vector. An "expression vector" is a construct that can be used to transform a selected host cell and provides for expression of a coding sequence in the selected host. Expression vectors can for instance be cloning vectors, binary vectors or integrating vectors. Expression comprises transcription of the nucleic acid molecule preferably into a translatable mRNA. Regulatory elements ensuring expression in prokaryotes and/or eukaryotic cells are well known to those skilled in the art. In the case of eukaryotic cells they comprise normally promoters ensuring initiation of transcription and optionally poly-A signals ensuring termination of transcription and stabilization of the transcript. Possible regulatory elements permitting expression in prokaryotic hosts comprise, e.g., the PL, lac, trp or tac promoter in *E. coli*, and examples of regulatory elements permitting expression in eukaryotic host cells are the AOX1 or GAL1 promoter in yeast or the CMV-, SV40-, RSV-promoter (Rous sarcoma virus), CMV-enhancer, SV40-enhancer or a globin intron in mammalian and other animal cells. In this

context, suitable expression vectors are known in the art such as Okayama-Berg cDNA expression vector pcDV1 (Pharmacia), pCDM8, pRc/CMV, pcDNA1, pcDNA3 (In-vitrogene), pSPORT1 (GIBCO BRL). Typical prokaryotic cloning and expression vectors include: plasmid vectors like the pUC series (e.g. pUC18, pUC19), pGEM series (e.g. pGEM 7zf+, Promega, USA), pET series (e.g. pET22B, Novagen, USA), pBBC 1MCS series, pNOF (GL Biotech Germany), pCR-TOPO series and pCR Blunt (Invitrogen, USA), pBluescript series, pCAL series and pBC series (Stratagene, USA); Fosmid vectors like pEpifos5 and pCC1 (Epicentre, USA); Cosmid vectors like the Expand series (Expand I, II, III) (Roche, Germany), SuperCos (Stratagene, USA), pOJ436; BAC vectors like pBeloBAC, pCC1BAC (Epicentre, USA).

An alternative expression system which could be used to express a cell cycle interacting protein is an insect system. In one such system, *Autographa californica* nuclear polyhedrosis virus (AcNPV) is used as a vector to express foreign genes in *Spodoptera frugiperda* cells or in *Trichoplusia* larvae. The coding sequence of a nucleic acid molecule of the invention may be cloned into a nonessential region of the virus, such as the polyhedrin gene, and placed under control of the polyhedrin promoter. Successful insertion of said coding sequence will render the polyhedrin gene inactive and produce recombinant virus lacking coat protein coat. The recombinant viruses are then used to infect *S. frugiperda* cells or *Trichoplusia* larvae in which the protein of the invention is expressed (Smith, J. Virol. 46 (1983), 584; Engelhard, Proc. Nat. Acad. Sci. USA 91 (1994), 3224-3227).

In plants, promoters commonly used are the polyubiquitin promoter, and the actin promoter for ubiquitous expression. The termination signals usually employed are from the Nopaline Synthase promoter or from the CAMV 35S promoter. A plant translational enhancer often used is the TMV omega sequences, the inclusion of an intron (Intron-1 from the Shrunken gene of maize, for example) has been shown to increase expression levels by up to 100-fold. (Mait, Transgenic Research 6 (1997), 143-156; Ni, Plant Journal 7 (1995), 661-676). Additional regulatory elements may include transcriptional as well as translational enhancers. Advantageously, the above-described vectors of the invention comprises a selectable and/or scorable marker. Selectable marker genes useful for the selection of transformed cells and, e.g., plant tissue and plants are well known to those skilled in the art and comprise, for example, antimetabolite resistance as the basis of selection for dhfr, which

confers resistance to methotrexate (Reiss, *Plant Physiol. (Life Sci. Adv.)* 13 (1994), 143-149); npt, which confers resistance to the aminoglycosides neomycin, kanamycin and paromycin (Herrera-Estrella, *EMBO J.* 2 (1983), 987-995) and hygromycin, which confers resistance to hygromycin (Marsh, *Gene* 32 (1984), 481-485). Additional selectable genes have been described, namely trpB, which allows cells to utilize indole in place of tryptophan; hisD, which allows cells to utilize histinol in place of histidine (Hartman, *Proc. Natl. Acad. Sci. USA* 85 (1988), 8047); mannose-6-phosphate isomerase which allows cells to utilize mannose (WO 94/20627) and ODC (ornithine decarboxylase) which confers resistance to the ornithine decarboxylase inhibitor, 2-(difluoromethyl)-DL-ornithine, DFMO (McConlogue, 1987, In: *Current Communications in Molecular Biology*, Cold Spring Harbor Laboratory ed.) or deaminase from *Aspergillus terreus* which confers resistance to Blasticidin S (Tamura, *Biosci. Biotechnol. Biochem.* 59 (1995), 2336-2338).

Useful scorable marker are also known to those skilled in the art and are commercially available. Advantageously, said marker is a gene encoding luciferase (Giacomin, *Pl. Sci.* 116 (1996), 59-72; Scikantha, *J. Bact.* 178 (1996), 121), green fluorescent protein (Gerdes, *FEBS Lett.* 389 (1996), 44-47),  $\beta$ -glucuronidase (Jefferson, *EMBO J.* 6 (1987), 3901-3907) or secreted alkaline phosphatase (SEAP) (Schlatter et al. (2001) *Biotechnol Bioeng.* 5, 75(5), 597-606). This embodiment is particularly useful for simple and rapid screening of cells, tissues and organisms containing a vector of the invention.

The present invention furthermore relates to host containing an aforementioned vector or an aforementioned nucleic acid molecule. Said host may be produced by introducing said vector or nucleotide sequence into the host by transfection or transformation wherein the nucleotide sequence and/or the encoded (poly)peptide/protein is foreign to the host. Upon the presence of said vector or nucleotide sequence in the host the expression of a protein encoded by the nucleotide sequence of the invention or comprising a nucleotide sequence or a vector according to the invention is mediated.

By "foreign" it is meant that the nucleotide sequence and/or the encoded (poly)peptide/protein is either heterologous with respect to the host, this means derived from a cell or organism with a different genomic background, or is

homologous with respect to the host but located in a different genomic environment than the naturally occurring counterpart of said nucleotide sequence. This means that, if the nucleotide sequence is homologous with respect to the host, it is not located in its natural location in the genome of said host, in particular it is surrounded by different genes. In this case the nucleotide sequence may be either under the control of its own promoter or under the control of a heterologous promoter. The vector or nucleotide sequence according to the invention which is present in the host may either be integrated into the genome of the host or it may be maintained in some form extrachromosomally. In this respect, it is also to be understood that the nucleotide sequence of the invention can be used to restore or create a mutant gene via homologous recombination.

Moreover, the present invention related to a method for producing a (poly)peptide as encoded by a nucleic acid molecule of the invention, comprising culturing the host of the invention under suitable conditions and isolating said (poly)peptide from the culture.

Isolation and purification of the recombinantly produced proteins and (poly)peptides may be carried out by conventional means including preparative chromatography and affinity and immunological separations involving affinity chromatography with monoclonal or polyclonal antibodies specifically interacting with said proteins/(poly)peptides. Preferably, said antibodies are antibodies of the invention as described herein below.

As used herein, the term „isolated protein“ includes proteins substantially free of other proteins, nucleic acids, lipids, carbohydrates or other materials with which it is naturally associated. Such proteins however not only comprise recombinantly produced proteins but include isolated naturally occurring proteins, synthetically produced proteins, or proteins produced by a combination of these methods. Means for preparing such proteins are well understood in the art. The proteins of the invention are preferably in a substantially purified form. A recombinantly produced version of said proteins, including secreted proteins, can be substantially purified by the one-step method described in Smith and Johnson, 1988.



An alternative embodiment of the invention relates to a (poly)peptide encoded by a nucleic acid molecule of the invention or as obtained by the method of the invention.

Preferably said (poly)peptide or fragment thereof is glycosylated, phosphorylated, amidated and/or myristylated.

Furthermore, the present invention relates to an antibody or an aptamer specifically recognizing the aforementioned (poly)peptide or a fragment or epitope thereof. Said antibody may be a monoclonal or a polyclonal antibody.

The term "fragment thereof" as used herein refers to fragments of said (poly)peptide/protein which are characterized by their capability to induce an immunological response in an immunized organism. Said response may be induced by the protein or fragment thereof either alone or in combination with a hapten, an adjuvant or other compounds known in the art to induce or elicit immunoresponses to a protein or fragment thereof.

The term "epitope" defines a single antigenic determinant. Said determinant is at least a portion of an antigen to which e.g. an antibody specifically binds to by its paratope; see Roitt et. al. (1993) Immunology 3<sup>rd</sup> edition, Mosby.

A preferred embodiment of the invention relates to an antibody which is a monoclonal antibody.

Said antibody, which is monoclonal antibody, polyclonal antibody, single chain antibody, or fragment thereof that specifically binds said peptide or polypeptide also including bispecific antibody, synthetic antibody, antibody fragment, such as Fab, a F(ab<sub>2</sub>)', Fv or scFv fragments etc., or a chemically modified derivative of any of these (all comprised by the term "antibody"). Monoclonal antibodies can be prepared, for example, by the techniques as originally described in Köhler and Milstein, Nature 256 (1975), 495, and Galfré, Meth. Enzymol. 73 (1981), 3, which comprise the fusion of mouse myeloma cells to spleen cells derived from immunized mammals with modifications developed by the art. Furthermore, antibodies or fragments thereof to the aforementioned peptides can be obtained by using methods which are described, e.g., in Harlow and Lane "Antibodies, A Laboratory Manual", CSH Press, Cold Spring Harbor, 1988. When derivatives of said antibodies are obtained by the

phage display technique, surface plasmon resonance as employed in the BIAcore system can be used to increase the efficiency of phage antibodies which bind to an epitope of the peptide or polypeptide of the invention (Schier, Human Antibodies Hybridomas 7 (1996), 97-105; Malmberg, J. Immunol. Methods 183 (1995), 7-13). The production of chimeric antibodies is described, for example, in WO89/09622. A further source of antibodies to be utilized in accordance with the present invention are so-called xenogenic antibodies. The general principle for the production of xenogenic antibodies such as human antibodies in mice is described in, e.g., WO 91/10741, WO 94/02602, WO 96/34096 and WO 96/33735. Antibodies to be employed in accordance with the invention or their corresponding immunoglobulin chain(s) can be further modified using conventional techniques known in the art, for example, by using amino acid deletion(s), insertion(s), substitution(s), addition(s), and/or recombination(s) and/or any other modification(s) known in the art either alone or in combination. Methods for introducing such modifications in the DNA sequence underlying the amino acid sequence of an immunoglobulin chain are well known to the person skilled in the art; see, e.g., Sambrook (1989), loc. cit..

The term "monoclonal" or "polyclonal antibody" (see Harlow and Lane, (1988), loc. cit.) also relates to derivatives of said antibodies which retain or essentially retain their binding specificity. Whereas particularly preferred embodiments of said derivatives are specified further herein below, other preferred derivatives of such antibodies are chimeric antibodies comprising, for example, a mouse or rat variable region and a human constant region.

The term "scFv fragment" (single-chain Fv fragment) is well understood in the art and preferred due to its small size and the possibility to recombinantly produce such fragments.

The term "specifically binds" in connection with the antibody used in accordance with the present invention means that the antibody etc. does not or essentially does not cross-react with (poly)peptides of similar structures. Cross-reactivity of a panel of antibodies etc. under investigation may be tested, for example, by assessing binding of said panel of antibodies etc. under conventional conditions (see, e.g., Harlow and Lane, (1988), loc. cit.) to the (poly)peptide of interest as well as to a number of more or less (structurally and/or functionally) closely related (poly)peptides. Only those antibodies that bind to the (poly)peptide/protein of interest but do not or do not

essentially bind to any of the other (poly)peptides which are preferably expressed by the same organism/tissue as the (poly)peptide of interest, e.g. by a crenarchaeote, are considered specific for the (poly)peptide/protein of interest and selected for further studies in accordance with the method of the invention.

In a further alternative embodiment the present invention relates to a transgenic non-human mammal whose somatic and germ cells comprise at least one gene encoding a functional polypeptide selected from the group consisting of:

- (a) the polypeptide of the invention;
- (b) a polypeptide having an amino acid sequence that is at least 60%, preferably at least 80%, especially at least 90%, advantageously at least 99% identical to the amino acid sequence of (a); and
- (c) a polypeptide having the amino acid sequence of (a) with at least one conservative amino acid substitution.

A method for the production of a transgenic non-human animal, for example transgenic mouse, comprises introduction of the aforementioned polynucleotide or targeting vector into a germ cell, an embryonic cell, stem cell or an egg or a cell derived therefrom. The non-human animal can be used in accordance with the invention in a method for identification of compounds, described herein below. Production of transgenic embryos and screening of those can be performed, e.g., as described by A. L. Joyner Ed., *Gene Targeting, A Practical Approach* (1993), Oxford University Press. The DNA of the embryonal membranes of embryos can be analyzed using, e.g., Southern blots with an appropriate probe; see supra. A general method for making transgenic non-human animals is described in the art, see for example WO 94/24274. For making transgenic non-human organisms (which include homologously targeted non-human animals), embryonal stem cells (ES cells) are preferred. Murine ES cells, such as AB-1 line grown on mitotically inactive SNL76/7 cell feeder layers (McMahon and Bradley, *Cell* 62:1073-1085 (1990)) essentially as described (Robertson, E. J. (1987) in *Teratocarcinomas and Embryonic Stem Cells: A Practical Approach*. E. J. Robertson, ed. (Oxford: IRL Press), p. 71-112) may be used for homologous gene targeting. Other suitable ES lines include, but are not limited to, the E14 line (Hooper et al., *Nature* 326:292-295 (1987)), the D3 line

(Doetschman et al., J. Embryol. Exp. Morph. 87:27-45 (1985)), the CCE line (Robertson et al., Nature 323:445-448 (1986)), the AK-7 line (Zhuang et al., Cell 77:875-884 (1994)). The success of generating a mouse line from ES cells bearing a specific targeted mutation depends on the pluripotency of the ES cells (i. e., their ability, once injected into a host developing embryo, such as a blastocyst or morula, to participate in embryogenesis and contribute to the germ cells of the resulting animal). The blastocysts containing the injected ES cells are allowed to develop in the uteri of pseudopregnant non-human females and are born, e.g. as chimeric mice. The resultant transgenic mice are chimeric for cells having either the recombinase or reporter loci and are backcrossed and screened for the presence of the correctly targeted transgene (s) by PCR or Southern blot analysis on tail biopsy DNA of offspring so as to identify transgenic mice heterozygous for either the recombinase or reporter locus/loci.

The transgenic non-human animals may, for example, be transgenic mice, rats, hamsters, dogs, monkeys, rabbits, pigs, or cows. Preferably, said transgenic non-human animal is a mouse.

The figures show:

Figure 1A shows the Structure of humic acids (after Stevenson 1982).

Figure 1B shows the Structure of fulvic acids (after Buffle 1977).

Figure 2 shows the Structure of polyvinylpyrrolidone (Monomer)

Figure 3 shows an example of an 2-phase PVP-low-melting agarosegel for simultaneous DNA purification and size resolution (Discontinuous Affinity-Gelelectrophoresis, DAG). The PVP-agarose-phase (the first phase) typically takes up approximately 1/4 to 1/3 of the gel but may take up maximally 80-95 % of the gel. After migrating through the PVP-phase into the agarose-only phase gel segments containing DNA fragments of interesting size may be excised and appropriately treated. The DNA migrates from minus (-) to plus (+).

Figure 4 shows an example of an 2-Phase Column chromatography for size resolution and affinity purification of DNA. In the first phase DNA in solution is resolved from inhibitors by size-exclusion chromatography. After passing this part of the composite column, the DNA passes a phase containing an affinity matrix (e.g. PVPP) to selectively bind and retard inhibitors. The DNA elutes from the column largest molecules first. The liquid flow is driven by hydrostatic or peristaltic pressure.

Figure 5 shows pulsed-field gelelectrophoretic (PFGE-) separations of metagenomic DNA isolated from soil (method A). In Figure 5 A the high molecular weight DNA is concentrated in a compression zone above 600 kbp. Yeast genomic DNA and a commercial size marker are added for reference. In Figure 5 B most DNA is sized to around 500 kbp but fragments up to 1.9 Mbp are visible.

M: Size marker (in kbp)

Figure 6 shows an insert analysis of metagenomic clones containing soil-derived DNA. Clones were digested with Not I.

Figure 7 shows an expression screening of arrayed fosmid clones containing metagenomic, soil-derived DNA. The encircled clone shows a halo of substrate degradation indicating hydrolase enzyme activity.

Figure 8 shows the quantification of metagenomic soil DNA in fractions eluting from columns after PVPP/sepharose 2B chromatography.

Figure 8 A: The DNA bands of fractions 6 - 22 were quantified by fluorescence intensity comparison after gel electrophoresis in a 1% agarose gel containing ethidium bromide. Gel documentation and analysis was done using GeneTools software from SynGene (UK). A 1% agarose gel was chosen to concentrate heterogenous DNA fragments in a single band for the sake of simplifying quantification. This was achieved yet at the price of size-resolution. The apparent comigration of eluting DNA with the 23 kbp marker band therefore is an underestimation of true maximum DNA fragment sizes. Lane 1:  $\lambda$  Hind III Marker DNA

Figure 8 B: Analysis of DNA in PVPP/sepharose 2B chromatography fractions by agarose gel electrophoresis

\* indicates fractions containing pure DNA

! indicates fractions containing large amounts of humic acids

Figure 9 shows the separation of soil metagenome DNA from humic substances by chromatography on PVPP/sepharose 2 B column.

Figure 9 A: Spectral absorption of eluting fractions 1–22 at 260 (DNA and humic acids) and 230 nm (humic acids) are plotted along their ratio and the relative DNA amounts as determined by agarose gel electrophoresis (see figure 8). A high A<sub>260</sub>/A<sub>230</sub> ratio indicates pure DNA with low humic/fulvic acid contamination as can be seen in fractions 8 and 9. Absorption at 260 nm and 230 nm rises in two peaks in later fractions (12 and 17) possibly due to two size-populations of humic acids.

Figure 9 B: The result of a corresponding experiment is shown in Figure 9 B. Microbial DNA from humic substances in crude extract from soil is separated by PVPP/sepharose 2B chromatography. The gel electrophoretic analysis of fractions derived from this separation as described for Figure 9 A

Figure 10: Restriction digest (*NotI*) of randomly selected environmental fosmid clones separated by pulse field gel electrophoresis. Lanes 1,2: DNA size standards. A band of 7.5 kb visible in all lanes corresponds to the fosmid vector.

Figure 11: Phylogenetic tree based on full-length 23S rDNA sequences of Bacteria, Archaea and of sequences obtained from marine environmental genomic clones. Different alignment filters were used to evaluate the phylogenetic reconstruction. The tree topology shown here is based on a maximum parsimony analysis (using 1048 conserved positions selected by a positional variability filter). Closed circles indicate branching points supported by different phylogenetic methods and filters in reconstructions with the 23S and the corresponding 16S rDNA tree (see Experimental Procedures in appended example 5).

Figure 12: Schematic representation of the archaeal fosmid clone 29i4/SEQ ID NO: 1. Different shadings indicate the phylogenetic affinity of the putative protein-coding genes to archaea (diag. stripes), bacteria (dots), bacteria and archaea (vertical stripes), or archaea, bacteria and eukarya (grey). Hypothetical genes with no homologs are shown without fillings. ORF numbers match those in Table 1.

Figure 13: Phylogenetic analysis with selected sequences of the ETF-like protein family. Homologs of FixA proteins from archaea form a monophyletic group with the FixA proteins of nitrogen fixing bacteria and of *Thermotoga maritima*, clearly distinguished from "housekeeping" ETF proteins (Etf $\beta$  homologs) of bacteria, archaea and eukaryotes. A third distinct subgroup (termed FixA paralog) is formed by a few as yet uncharacterized sequences from bacteria and archaea. For details of the phylogenetic reconstruction see Experimental Procedures in appended example 5.

Figure 14:

Figure 14 A: Size selection of soil DNA in a standard agarose gel

Preparative standard PFGE gel (EtBr stained) separation of soil DNA. The compression zone above 100 kbp was excised and resulting agarose plugs subsequently subjected to enzymatic digestion to test purity and clonability (Figure 14 B).

Figure 14 B: Failed hydrolyses of soil DNA (*Hind*III) purified in a standard agarose gel

Analytical PFGE showing DNA from an enzyme titration experiment (EtBr stained). DNA containing cubic plugs cut out from a standard PFGE agarose gel (Figure 14 A) and containing fragments sized over 100 kb were added to a reaction mix containing increasing units of *Hind*III restriction enzyme. After terminating the reaction the DNA was re-run on a second PFGE. It is evident from the similarly undigested DNA at 0U and at 100U enzyme concentration that this DNA is highly resistant to digestion with *Hind* III indicating inhibiting impurities that prevent molecular manipulation and cloning.

Figure 14 C: Size selection and purification of soil DNA in a 2-phase gel

Preparative 2-phase PFGE gel (SYBR green stained, inverted image) purification of soil DNA. The compression zone of DNA above 250 kbp was excised and subsequently subjected to enzymatic digestion as a proof of purity and clonability (Figure 14 D).

Figure 14 D: Successful partial hydrolyses of soil DNA (*Hind*III) purified in a 2-phase gel

Analytical PFGE showing soil DNA from an enzyme titration experiment (EtBr stained). 2-phase gel purified soil DNA cut out from a preparative 2-Phase Gel (Figure 14 C) and sized over 280 kb was *Hind* III digested and re-run on a second PFGE. It is evident that after purification in the 2-phase gel the DNA is readily digested to smaller sizes even with as little as 1U *Hind* III indicating a high degree of purity and access by enzymes necessary for subsequent molecular cloning.

## Examples

The following examples illustrate the invention. These examples should not be construed as limiting; the examples are included for purposes of illustration and the present invention is limited only by the claims.

### Example 1 Generation of a fosmid library from soil metagenomic DNA

#### 1.1 Preparation of High Molecular Weight DNA from Soil (A)

Soil was collected from the upper layer of a partially ruderalized sandy ecosystem in Weiterstadt (Germany). About 50 g were suspended in 300 ml of buffer (pH 8, 20 mM Tris-HCl, 10 mM  $\epsilon$ -aminocaproic acid, 10 mM EDTA) and incubated at 4°C for 15 h with gentle shaking. The suspension was sieved to remove larger particles and after centrifugation of the filtrate (5000 x g 30 min.) the resulting microbial fraction was embedded in agarose plugs (0.5 % low-melting SEAPlaque, FMC Bioproducts). The plugs were incubated at 37°C for 1 hour in lysozyme buffer (100 mM EDTA pH 8.0, 10 mM Tris-HCl pH 8.0, 50 mM NaCl, 0.2% Na-deoxycholate, 1% laurylsarcosin, 2 mg/ml lysozyme), then transferred into ESP solution containing 2 mg/ml proteinase K, 1% lauroyl Sarcosin and 0.5 M EDTA and incubated at 50°C for 24 h under gentle rotation and with 1 exchange of ESP buffer. 3 agarose plugs were placed in a 1%



agarose gel (Sigma A-2929) that contained 2 % PVP (Sigma PVP-360 ) in the upper part and no PVP in the lower part. Pulsed-field electrophoresis was performed in a CHEF-DR II PFGE machine (BioRad) at 10°C, 6 V cm<sup>-1</sup> for 20 h with 1 to 4 s pulses (Figure 5). A slice of agarose containing DNA in the size range of >30 kbp was cut out of the gel inserted into appropriately sized slots in a second gel and re-electrophoresed for a second size selection. After electroelution the resulting DNA was dialyzed and concentrated in a microconcentrator (Vivascience).

## **1.2 Preparation of High Molecular Weight DNA from Soil (B)**

Alternatively total soil DNA was extracted using a protocol modified from Zhou and coworkers (Zhou et al., (1996) loc. cit.). Soil (5 g) was resuspended in 13.5 ml extractionbuffer (100 mM Tris/HCl pH8.0, 100 mM Na-EDTA, 100 mM Na-phosphate pH 8.0, 1.5 M NaCl, 1% CTAB) followed by an optional 3 cycles of freezing in liquid nitrogen and boiling in a microwave oven. After adding 1.5 ml of lysozyme solution (50 mg/ml) and 30 min. incubation at 37°C, 200µl proteinase K solution (10 mg/ml) were added followed by another 30 min. incubation at 37°C. Then 3 ml 10 % SDS was added followed by 2 hours incubation at 65°C. After centrifugation (10 min., 6000 x g, room temperature) the supernatant was collected and the pellet reextracted twice for 10 min. at 65°C with 4.5 ml of extraction buffer and 1 ml of 10% SDS. All supernatants were united and extracted with an equal volume of chloroform/isoamylalcohol (24:1 vol/vol). DNA was precipitated from the aqueous phase with 0.6 vol isopropanol, pelleted by centrifugation (16000 x g, 20 min., room temperature), washed in 70% ethanol and dissolved in 200µl TE buffer (10 mM Tris-HCl pH 8.0, 1 mM EDTA). Depending on the soil sample this solution was yellowish to dark brown. This DNA routinely defied any enzymatic manipulation by restriction enzymes and could be successfully used for PCR only at dilutions below 1/1000. It was therefore electrophoresed either in a constant voltage electrophoresis using a 0.5 % agarose 2-phase gel with 2 % PVP in the first phase (4 hours at 60 V ) or in a 2-phase PFGE gel as described in (A).

**Example 2: Library construction and analysis**

About 0.5 µg of purified DNA was enzymatically treated to prepare 5' phosphorylated blunt ends and was ligated to the linearized and dephosphorylated fosmid vector pEpiFOS-5 (pEpiFOS Fosmid Library Production Kit, Epicentre). After in vitro packaging into lambda phages (Epicentre) and infection of E.coli strain EPI100 (Epicentre), cells were plated on LB medium containing 12.5 µg/ml chloramphenicol. The colonies were transferred to individual wells of 384-well microtitre plates containing 50 µl of LB with 12.5 µg/ml chloramphenicol and 7 % glycerol (v/v) and were incubated at 37°C for 24 hours. The library was stored at -80°C.

After blunt-end cloning into the fosmid vector pEpiFOS-5 about 50000 colonies were obtained per µg of soil DNA. Our final library was constructed from a single ligation mixture and contained 25278 clones arrayed in 66 384-well microtitre plates.

Restriction analysis of 30 randomly chosen clones with NotI showed insert sizes between 32.5 and 43.5 kbp, with an average of 36.5 kbp which corresponds well with the insert size range acceptable for this type of vector (Figure 6). The total library therefore contained an estimated 0.9 Gbp of environmental genomic soil DNA, which represents 225 genome equivalents assuming a 4 Mbp average genome size. Most inserts analyzed exhibited complex patterns after NotI digestion, suggesting that these clones contained DNA with high GC content.

**Example 3: Screening of the library**

The arrayed clones of the library were plated

- a) onto LB-agar to grow the cells for a subsequent preparation of fosmid-pools as a resource for sequence-based screenings (using PCR and degenerate primers to generate metagenome sequence tags e.g. as probes for hybridisation) and
- b) onto LB-agar containing specific substrates for the detection of enzymatically active colonies e.g. through scoring of clearing zones around the colonies (Figure 7) and
- c) onto LB-agar for growth and subsequent overlay with a lawn of indicator organisms to score for recombinants producing antibiotics.

**Example 4: Purification of soil metagenomic DNA by use of a 2-phase gel-permeation/affinity column**

Metagenome DNA from soil was extracted by gentle chemical lysis as described before (Zhou et al., (1996) *loc. cit.*). Crude DNA extract was passed over a 2-phase column by gravity-flow using a borosilicate glass column (BIORAD #737-0717). The column (7x150 mm) was packed with a lower phase of 1ml PVPP (SIGMA # P-6755) and an upper layer of Sepharose 2B (SIGMA # 2B-300) to a final volume of 5 ml. After equilibrating the column with 20ml of running buffer (100 mM NaCl, 10 mM Tris, 1 mM EDTA; pH 8.0) separation of DNA from humic and fulvic acids was initiated by applying crude DNA extract (ideally 1-5 % resin volume) to the top of the column. Eluting fractions of 300 µl were collected dropwise (6 drops) and subsequently analysed. Relative DNA content was quantified by comparing the fluorescence of DNA bands in each fraction after electrophoresis on a 1% agarose gel containing ethidium bromide (figure 8) and documentation using a CCD camera system (GeneGenius) and GeneTools software package (both Syngene, UK). The spectral absorption of fractions at 230 nm and 260 nm was measured spectrophotometrically. Fractions 8 and 9 contain pure clonable DNA. Analysis of later eluting fractions 10–12 show humic acid contamination as can be seen in the decrease in the ratio OD260/OD230. This is due to the significant rise of the OD230 values - the maximum absorption for humic acids. As humic substances also absorb at longer wavelengths the OD260 values increase similarly. The separation of fractions containing pure DNA and fractions containing humic substances can be improved further by adjusting the ratio of the loaded volume of crude DNA extract to the resin volume of the column in a way known to those skilled in the art and by adjusting the DNA content of the loaded samples.

**Example 5: Exemplified identification of genes of a microorganism as part of a metagenome in accordance with the invention.**

Traditionally, soil microbiology has focused on the description of cultivable microorganisms, while functional aspects of soil microbial communities have mainly been restricted to bulk studies that involved the monitoring of substrates and levels of end products. The application of molecular techniques to microbial ecology

revealed that many of the microbial transformations in the environments might be performed by organisms that have not yet been cultivated and thus far remained uncharacterized (Pace (1997) *Science* 276, 734-740, Hugenholtz et. al. (1998) *J. Bacteriol.* 180:4765-74). Soil was confirmed to be particularly rich in microbial diversity based on phylogenetic studies with 16S rRNA genes directly amplified from environmental samples (see e.g. Hugenholtz et. al. (1998) loc. cit., Borneman et al. (1996) *Appl. Environm. Microbiol.* 62, 1935-1943., Barns et. al. (1999) *Appl. Environm. Microbiol.* 65, 1731-1737, Dunbar et. al. (1999) *Appl. Environm. Microbiol.* 65, 1662-1669). Frequently, evidence was even found for the existence of microorganisms from divisions that were not previously associated with soil habitats and of which no cultivated relatives. Among one of the most striking discoveries was the frequent detection of non-thermophilic members of the archaeal kingdom Crenarchaeota (DeLong (1998) *Curr. Opin. Genet. Dev.* 8, 649-654). 16S rDNA sequences of these archaea were first identified in marine picoplankton (DeLong (1992) *Proc. Natl. Acad. Sci. U S A* 89, 5685-9, Fuhrman et. al. (1992) *Nature* 356, 148-9) and then found in freshwater habitats (Hershberger et. al. (1996) *Nature* 384, 420, Schleper et. al. (1997a) *Appl. Env. Microbiol.* 63, 321-323, McGregor et. al. (1997) *Appl. Env. Microbiol.* 63, 1178-1181, Jurgens et. al. (2000) *FEMS Microbiol. Ecol.* 34, 45-56) and in soils from various locations in the United States (Bintrim, et. al. (1997) *Proc. Natl. Acad. Sci. USA* 94, 277-282, Buckley et. al. (1998) *Appl. Environ. Microbiol.* 64, 4333-4339), Finland (Jurgens et. al. (1997) *Appl. Environ. Microbiol.* 63, 803-805), Japan (Ueda et. al. (1995) *Eur. J. Soil Sci.* 46, 415-421, Kudo et. al. (1997) *Biosci. Biotechnol. Biochem.* 61, 917-20) and Germany (Sandaa, et. al. (1999) *Appl. Environ. Microbiol.* 65, 3293-3297, Ochsenreiter et. al., in preparation). The ubiquitous ecological distribution of crenarchaeota was very surprising, because their cultivated relatives are exclusively (hyper)thermophiles isolated from terrestrial and marine hot springs.

Quantitative estimates have demonstrated the significant occurrence of non-thermophilic crenarchaeota in marine habitats (Massana et. al. (1997) *Appl. Environ. Microbiol.* 63, 50-6, Karner et. al. (2001) *Nature* 409, 507-510), in freshwater sediments (McGregor et. al. (1997) loc. cit.) and in soil (14,19). Some crenarchaeotal lineages were shown to be specifically associated with plant roots, indicating that the

organisms might play a role in the ecology of the rhizosphere (Simon et. al. (2000) *Environ. Microbiol.* 2, 495-505, Chelius & Triplett (2001) *Microb. Ecol.* 41,252-263). While 16S rRNA studies have provided insights into the huge extent of microbial diversity novel approaches are being sought to be able to functionally characterize those microorganisms that have escaped classical cultivation approaches.

Inspired by the rapid advances in microbial genomics of cultivated organisms, a novel approach has recently been initiated to characterize uncultivated organisms that have solely been predicted in rRNA gene surveys. It involves the construction of complex habitat-specific gene libraries by direct cloning of genomic fragments from environmental samples into cloning vectors (DeLong (2001) *Curr. Opin. Microbiol.* 4, 290-295). With the help of phylogenetically relevant gene markers, such as e.g. rDNA genes, large genomic fragments of specific phylotypes can be isolated from these libraries. Yet at present the full potential of this approach cannot be realized as technical constraints severely hamper the cloning of large DNA inserts, particularly from microbial consortia of inhibitor-rich environments like soils or sediments. The approach has successfully been applied to characterize uncultivated, marine microorganisms: Several genome fragments of the symbiotic crenarchaeote *Cenarchaeum symbiosum* and of marine archaea representing abundant components of the picoplankton in North Pacific and Antarctic waters were retrieved from BAC and fosmid libraries (Stein et. al. (1996) *J. Bacteriol.* 178,591-599, Schleper et. al. (1998) *J. Bacteriol.* 180, 5003-5009, Béjà et. al. (2000a) *Environ. Microbiol.* 2, 516-529, Beja et. al. (2002a) *Nature* 415:630-633, Béjà et. al. (2002b) *Appl. Environ. Microbiol.* 68, 335-45). A comparison of crenarchaeotal fosmids revealed significant genomic divergence even in clones with identical 16S rRNA sequences (Béjà et. al. (2002b) loc. cit.). The diversity of large photosynthetic gene clusters of proteobacteria was analyzed from marine planktonic genomic samples (Beja et. al. (2002a) loc. cit.). A novel type of rhodopsin, termed proteorhodopsin that functions as a light-driven proton pump was discovered in the genomic fragment of an uncultivated marine  $\gamma$ -proteobacterium (Béjà et. al. (2000b) *Science* 289, 1902-1906). The analyses of large genomic regions of hitherto uncultivated organisms also provided the basis for functional studies, including the monitoring of protein activities in the environment (Beja et. al. (2002a) loc. cit., Beja, et. al. (2001) *Nature* 411, 786-

9) and the characterization of proteins after expression in the surrogate host *E.coli* (Béjà et. al. (2000b) loc. cit., Schleper et. al. (1997b) *J. Bacteriol.* 179, 7803-7811).

The colocalization of functional, metabolic genes with phylogenetically ascribable genetic markers like rRNA genes provides insights into the physiological potential of uncultivated microorganisms. Clearly the likelihood of a physical colocalization of such markers on a contiguous cloned DNA stretch will be directly linked to DNA fragment size. This highlights the relevance of cloning large uninterrupted DNA fragments which is technically very difficult to achieve particularly from microbial consortia of inhibitor-rich environments like soils or sediments. We have developed procedures for the efficient purification of large DNA fragments by eliminating the polyphenolic compounds that heavily contaminate soil DNA. We have constructed complex genomic libraries and used these to isolate fragments from non-thermophilic crenarchaeota. While direct cloning of large DNA from soil samples has been demonstrated earlier (Rondon et. al. (2000) *Appl. Environ. Microbiol.* 66, 2541-2547), our study represents the first genomic characterization of a lineage of soil microorganisms that has solely been predicted by PCR-based studies.

## 5.1 Experimental Procedures

### *Preparation of DNA from Soil*

Soil was collected from the upper layer (0 to 5 cm) of a partially ruderalized sandy ecosystem ("Am Rotböhl") near Darmstadt (Germany) in early Spring 2001. DNA was prepared as described in Example 1.1, supra. About 50 g were suspended in 300 ml of buffer (20 mM Tris pH 8, 10 mM  $\epsilon$ -aminocaproic acid, 10 mM EDTA ) and incubated at 4°C for 15h with gentle shaking. The sample was filtered to remove larger particles, and the microbial fraction was centrifuged and embedded into agarose plugs (0.5% low-melting SEAPlaque, FMC Bioproducts). These plugs were incubated at 37°C for 1 hour in 100 mM EDTA, 10 mM Tris-HCl pH 8, 50 mM NaCl, 0,2% deoxycholate, 1% lauroyl sarcosine, 1 mg/ml lysozyme, then transferred into ESP buffer (2 mg/ml proteinase K, 1% lauroyl sarcosine, 0.5 M EDTA) and incubated at 50°C for 24 h with gentle rotation and with one exchange of buffer. Agarose plugs were placed in a 1% agarose gel (Sigma A-2929) which contained 2% polyvinylpyrrolidone (VP-360, Sigma) in the first half and no PVP in the second half. Pulse field gel electrophoresis was performed at 10°C, 6 V cm<sup>-1</sup> for 20 h with 1 to 4

sec pulses in a CHEF-DR II (BioRad). DNA of > 30 kb was extracted from the gel and submitted to a second size selection using a regular agarose gel. After electroelution the resulting DNA was dialyzed and concentrated in a microconcentrator.

#### *Library construction*

Purified DNA (0.5 µg) was enzymatically treated to prepare 5' phosphorylated blunt ends and was ligated into fosmid vector pEpiFOS-5 (pEpiFOS™, Epicentre). After *in vitro* packaging into lambda phages, the infected cells were plated on LB<sup>+</sup> medium (containing 12.5 µg/ml chloramphenicol). The colonies were transferred to 384-well microtitre plates containing 50 µl of LB<sup>+</sup> medium and 7 % glycerol (v/v). The plates were incubated at 37°C for 24 hours.

#### *16S rDNA diversity studies*

Primers specific for the domain Archaea (20F/958R, DeLong (1992) *Proc Natl Acad Sci USA* 89: 5685-9) and Bacteria (27F/1391R, Reysenbach and Pace (1995) *Archaea: a laboratory manual* (Cold spring Harbor Laboratory Press)) were used to amplify 16S rDNA fragments from the DNA used for constructing the large-insert library. The fragments were subsequently cloned into pGEM-T-easy (Promega) and sequenced. The ARB-software package (Ludwig et al. (1998) *Electrophoresis* 19: 554-568.) was used for alignments and phylogenetic analyses of the partial 16S rDNA genes.

#### *Screening and sequence analysis of Fosmid clone 29i4*

Plasmid DNA from the library was prepared from pools of 384 clones and screened by PCR with archaea-specific 16S rDNA primers (DeLong (1992) loc. cit.). A product of correct size (950 bp) was obtained from pool 29. It was randomly labelled with digoxigenin (Roche Biochemicals) and used as a probe in colony hybridization to identify the individual clone (i4). A subclone library was prepared from the fosmid DNA by mechanical shearing and cloning of 2-3 kbp fragments into pGEM-T-easy (Promega). The ends of the cloned DNA fragments were sequenced with vector primers using ABI3700 capillary sequencers. Remaining sequence gaps were closed by primer walking with sequence-derived oligonucleotides.

#### *Sequence annotation*

The ORF identification and automatic gene annotation were done with the help of the MAGPIE program package (Gaasterland and Sensen (1996) *Biochimie* 78: 302-310). The Wisconsin Package (Heidelberg Unix Sequence Analysis Server, HUSAR) was used for additional searches with GCGBLAST, identification of PFAM domains and transmembrane segments, for secondary structure prediction and peptide motifs. A tRNA gene was identified using the tRNA scan server (<http://www.genetics.wustl.edu/eddy/tRNAscan-SE/>). Multiple alignments were done with PILEUP and CLUSTAL and manually corrected in SEQLAB. The sequence was deposited in EMBL under the accession no. AJ496176.

### *Phylogenetic Analyses*

The ARB-software package (Ludwig et al. (1998) loc. cit.) was used for alignments and phylogenetic analyses of full-length 16S and 23S rDNA genes from Archaea and from the marine environmental clones of a euryarchaeote 37F11 (Béjà et al. (2000a) loc. cit.), and of two crenarchaeota (4B7, Stein et al. (1996) loc. cit., Béjà et al. (2002b) loc. cit.). The topologies of the 23S rDNA tree were evaluated using the maximum parsimony (parsimony interactive) and the distance matrix (Felsenstein correction) method with different alignment filters (gap-filter, positional variability filter, maximum frequency filter). The topologies of the corresponding 16S rDNA tree were evaluated using the same methods and in addition the maximum likelihood (fastDNAmI) method with different alignment filters as described above. Phylogenetic analysis of the putative FixA gene (ORF12/SEQ ID NO: 24) was performed using the protein parsimony (PROTPARS) and neighbor-joining (NEIGHBOR) programs from PHYLIP version 3.6 and PAUP version 3.1.1. The same overall topology was found with both methods based on 211 conserved positions from a sequence alignment of 33 FixA/ETFB (flavoprotein containing electron transport chain) homologs.

## **5.2 Results**

### *Construction of a fosmid library from soil DNA*

For preparation of high molecular weight DNA biomass from a soil sample was embedded into agarose plugs prior to lysis. The resulting DNA appeared heavily contaminated with polyphenolic compounds (i.e. humic and fulvic acids) as indicated by a dark brownish appearance. In the process of optimizing the subsequent



purification steps, a pulse field electrophoresis procedure was developed that involved an agarose gel with two phases. It allowed purification of the DNA from soil substances through a PVP (polyvinylpyrrolidone) containing phase. In a second phase the PVP was subsequently eliminated, while a first size selection of the DNA was achieved. Highly concentrated, pure and clonable DNA in the size range of 30 to 100 kb was recovered in this one-step electrophoresis procedure, thereby minimizing shearing effects that tend to occur in repeated electrophoresis procedures. The approach has been successfully applied to different soil samples, like ruderal, agricultural and forest soils, for rapid preparation of pure and concentrated high molecular weight DNA (data not shown). After blunt-end cloning into the fosmid vector about 50,000 colonies were obtained per  $\mu\text{g}$  of soil DNA. Our final library contained 25,278 clones.

Restriction analysis of 30 randomly chosen clones using *NotI* showed insert sizes between 32.5 and 43.5 kb. The library therefore contained approximately 0.9 Gbp of environmental genomic soil DNA. 27 of 30 inserts analyzed exhibited complex patterns after *NotI* digestion, suggesting that these clones contained DNA with high G+C content (*NotI* recognition sequence: GCGGCCGC; Fig. 10).

Sequencing of insert ends from 2688 clones revealed in about 25 % of the sequences significant similarities to protein genes from the data bases (e-values of  $< 10^{-10}$  in blastx searches). Among these were homologs of proteins from lineages that are typically found in soils, i.e. streptomycetes, clostridia and bacilli (data not shown). However most of the sequences did not show significant similarities to known protein genes. Together, these results confirmed that a great diversity of genomic DNAs was contained in the library.

To further monitor the diversity of the DNA used for construction of the library, a PCR-based 16S rDNA survey was performed using a primer set specific for Bacteria. 16S rRNA gene fragments affiliated with eight different bacterial phyla were identified in a random sample of 50 different sequences, many of which are typical for soil microbial assemblages, i.e. Actinobacteria, Chloroflexi,  $\alpha, \beta, \epsilon$  Proteobacteria, Planctomycetes, Acidobacterium/Holophaga, Cytophaga.

*Identification and analysis of a genomic clone from non-thermophilic crenarchaeota in soil.*

Using a multiplex PCR approach and 16S rDNA-specific probes, an archaeal fosmid clone was identified in the library. The insert of clone 29i4/SEQ ID NO: 1 was entirely sequenced and comprised 33,925 bp with an average G+C content of 40 %. It encoded a complete 16S and 23S ribosomal RNA operon, one tRNA gene and 17 predicted protein-encoding genes. The 16S RNA gene was 97 % identical over 711 positions (*E.coli* positions 8-719) to sequences previously recovered in a PCR study from the same soil (Ochsenreiter and Schleper, manuscript in preparation) and 95-97 % identical to sequences obtained from a soil in Wisconsin (Bintrim et al. (1997) loc. cit.). The ribosomal RNA operon of clone 29i4 consisted only of the 16S and 23S rRNA genes, without linked 5S rRNA or tRNA genes similar to cultivated thermophilic and uncultivated marine crenarchaeota. Phylogenetic analyses with the complete 16S rRNA and 23S rRNA genes confirmed the affiliation of clone 29i4 with the crenarchaeota (Fig. 11). As predicted by phylogenetic analyses with partial 16S rRNA sequences, it formed a sister group to the uncultivated marine organisms. The phylogenetic tree in Fig. 11 is based on complete 23S rRNA genes from known archaeal genomes and from environmental genomic fragments of marine archaea. The same branching orders were found in 16S rRNA phylogenies (see legend to Fig 11).

Ten of the 17 predicted protein-encoding genes showed significant similarity to genes of known function, two were conserved hypothetical genes, five open reading frames did not show any similarity to sequences in the databases (Febr. 2002). Eight of the predicted proteins showed highest similarities to archaeal homologs (Table 1 and Fig. 12). A family B DNA polymerase shared 46 % identical positions with its closest homolog from *Cenarchaeum symbiosum* (Schleper et al. (1997) J Bacteriol. 179, 7803-11). Although the C-terminal end of about 90 amino acids was not encoded on fosmid 29i4, all conserved exonuclease and polymerase motifs typically found in this class of DNA polymerases were identified in the deduced amino acid sequence.

Two other predicted proteins belonging to conserved archaeal protein families were an asparagine synthetase and a phosphoserine phosphatase, both involved in amino acid metabolism. Two putative glycosyl transferases (ORF07/SEQ ID NO.: 14 and 10/SEQ ID NO.: 20) shared significant similarities with homologs from the crenarchaeote *Sulfolobus solfataricus* and from the euryarchaeote *Pyrococcus* ssp.,

repectively. In contrast, a putative polyhydroxyalkanoate synthase (ORF04/SEQ ID NO.: 8) and a second  $\alpha/\beta$ -hydrolase (ORF02/SEQ ID NO.: 4) were most closely related to bacterial proteins.

A gene cluster was identified with high similarity in structure and sequence to the fixABCX operons found in many symbiotic nitrogen-fixing soil bacteria. Based on its similarity to the components of the flavoprotein-containing electron transport chain (ETF) that is involved in  $\beta$ -oxidation of fatty acids in mitochondria and some bacteria, the operon was proposed to encode a flavoprotein-containing electron transport chain (Weidenhaupt et al. (1996) *Arch Microbiol* 165: 169-178). In symbiotic bacteria the operon is co-regulated with other genes involved in nitrogen fixation (Gubler and Hennecke (1998) *J Bacteriol* 170: 1205-1214). FixABCX genes have also been identified in the genomes of several other bacteria and some thermophilic and hyperthermophilic archaea, i.e. in *Sulfolobus solfataricus*, *Thermoplasma acidophilum*, *Pyrobaculum aerophilum* and *Aeropyrum pernix*. Phylogenetic analyses of the putative FixA gene from 29i4 and homologs from completely sequenced microbial genomes indicated a close affiliation of FixA from 29i4 with other archaeal proteins. Together with FixA proteins from nitrogen fixing bacteria they formed a distinct subgroup within the Etf $\beta$ /FixA superfamily (Fig. 13).

A sensory histidine kinase was identified in close proximity but oriented in opposite direction to the fixABCX operon on clone 29i4. While the C-terminal half of the protein exhibited the conserved motifs typically found in sensory histidine kinases, no similarities to known proteins that might give hints to its specific role in sensing were found in the 350 amino-acid long N-terminal part.

Using PCR primers targeting the ends of the insert of clone 29i4 and internal protein coding genes, no contiguous genomic fragments overlapping clone 29i4 could thus far be detected in the library. However, additional archaeal clones were identified from non-thermophilic crenarchaeota in a second library that contained another 1.5 Gbp of DNA from the same soil sample. One of these clones was identified with archaea-specific 16S rRNA probes (as used for clone 29i4) and two other clones were identified through sequencing of insert ends from 768 randomly chosen clones. The sequence analysis of these genomic fragments is under way.

## 5.2 Discussion

The direct cloning of high-molecular weight DNA from soil is particularly difficult due to the occurrence of polyphenolic compounds that co-purify with DNA and severely inhibit PCR amplification reactions, hydrolysis by restriction enzymes, ligation and cloning procedures (Trevors and Van Elsas (1995) *Nucleic Acids in the Environment: Methods and Applications*, Springer Verlag, Berlin; Young et al. (1993) loc. cit.). Therefore, any purification protocol for DNA from soil samples must remove the phenolic compounds. Protocols have been developed that involve the addition of hexadecylmethylammoniumbromide (CTAB, Zhou et al. (1996) loc. cit.) or polyvinylpyrrolidone (PVP, Trevors and Van Elsas (1995) loc. cit.) in the extraction buffers, because these compounds complex polyphenolics or reduce their electrophoretic mobility when included in electrophoresis procedures (Young et al. (1993) loc. cit.). However, such compounds in turn have to be efficiently eliminated (e.g. by extraction, electrophoresis or affinity chromatography) because they inhibit enzymatic treatments of the isolated nucleic acids. Due to these difficulties most purification procedures either result in high quality DNA suitable for PCR amplification of gene fragments but too highly degraded for cloning of large fragments or it results in high molecular weight DNA that is not pure enough for cloning procedures. (refs above and own observations). Therefore, the device of the invention was developed for the electrophoresis procedure described herein above (two phases in which the DNA is first purified from polyphenolics in a PVP containing phase and subsequently cleaned in a second phase, thereby minimizing shearing effects that occur in repeated electrophoresis procedures). This technique allows reproducible obtainment of highly concentrated and pure, high molecular weight DNA. While purification techniques similar to those described by Zhou et al. (1996) loc. cit.) that was used by Rondon et al. (2000) loc. cit.) for cloning large DNA fragments were not applicable to many of our samples, the novel PVP electrophoresis procedure yielded reliable results with different soil samples of varying organic and humic acid content. Successful purification of the DNA was independent of the lysis procedures that we used, i.e. direct lysis of soil samples as in Zhou et al. (1996) loc. cit.) or lysis of microbial fractions as described here.

The complex environmental libraries constructed exemplarily by using the device and method of the invention contain a large fraction of the total genomic content of a soil microbial population, which has been referred to as the soil "metagenome" (Rondon et al. (2000) loc. cit.). The library characterized here was constructed in a BAC-derived fosmid containing cos-sites for *in vitro* packaging with lambda phages. Said vector yielded significantly larger clone numbers than using classical BAC-vectors (data not shown) and it allowed the direct cloning of undigested DNA by blunt-end ligation, thereby avoiding any bias introduced by restriction digests. Using 50 g of soil, a library with 0.9 Gbp of environmental DNA was constructed which represents approximately 225 genome equivalents assuming a 4 Mbp average genome size. Using archaea-specific 16S rDNA probes the fosmid clone 29i4 was identified. Sequence analysis demonstrated that a contiguous genomic fragment of non-thermophilic soil crenarchaeota was isolated: (i) Phylogenetic analyses based on the complete ribosomal 16S and 23S RNA genes indicated the specific affiliation with the crenarchaeotal clade as predicted in PCR-based studies. (ii) genes affiliated with archaea were found dispersed over the entire clone insert (iii) G+C content and codon usage of the predicted protein genes were similar over the entire insert and (iv) the deduced aminoacid sequence of a DNA polymerase gene showed greatest similarity to its homolog from the uncultivated marine symbiont *Cenarchaeum symbiosum*. Functional and biochemical analysis of the latter protein had confirmed the predicted non-thermophilic phenotype of this crenarchaeote (Schleper et al. 1997) J Bacteriol. 179, 7803-11).

On the other hand, significant differences to the content and structure of genomic fragments from uncultivated non-thermophilic marine archaea revealed that crenarchaeota from soil have significantly diverged from their relatives in other environments. An unusually large gap in the 16S-23S RNA operon and the lack of a GSAT gene (glutamate semialdehyde aminotransferase), which was consistently found to be directly linked to the ribosomal operon on all marine crenarchaeotal genome fragments (Béjà et al. (2002) loc. cit.) indicates the difference. Direct comparison of the soil fosmid clone to the genomic clones obtained from marine planctonic crenarchaeota and from the symbiont *C. symbiosum* revealed only one related protein encoding gene, i.e. the putative DNA polymerase. Genes on clone 29i4/SEQ ID NO: 1 appeared to be less densely packed than in genomes of other

archaea. Only 69 % of the sequence encoded RNA or protein genes. There was the large intergenic region in the 16S/23S rRNA cluster of 830 bp, which is atypical for ribosomal operons in crenarchaeota. Another large non-coding region of 2787 bp was found between ORF10/SEQ ID NO: 20 and fixA. No apparent genes or distinctive structural features, e.g. repetitive elements were identified in the non-coding regions.

The genomic information contained on fosmid 29i4/SEQ ID NO: 1 gives first insights into metabolic properties of crenarchaeota from soil and can serve as a basis for functional genomic studies. Beside genes encoding proteins for "house-keeping" functions (replication, aminoacid metabolism), two  $\alpha/\beta$  type hydrolases so far seem to be particularly found in soil crenarchaeota. One of them encodes a putative protein involved in the synthesis of polyhydroxyalkanoates. The operon encoding FixABCX revealed a putative flavoprotein containing electron transport chain that is commonly found in symbiotic nitrogen-fixing bacteria. A detailed phylogenetic analysis indicated that the putative FixA protein of 29i4 is most closely affiliated with archaeal FixA homologs and not with the FixA proteins of Bacteria or the paralogous ETF proteins from other species of Archaea, Eucarya or Bacteria. Therefore, it seems unlikely that the fixABCX genes of 29i4 have been acquired by horizontal gene transfer e.g. from symbiotic nitrogen-fixing bacteria which might reside in the same soil habitat. They rather seem to originate from a common ancestor of the Archaea. None of the obligately aerobic archaea that contain the fixABCX genes is known to be capable of fixing nitrogen. Expression analysis of this operon in well-studied thermophilic model organisms, such as *Sulfolobus solfataricus* might shed light on its physiological role in crenarchaeota.

#### **Example 6: Separation of DNA extracts on a PVPP/Sepharose 2B Column**

Metagenome DNA from soil was extracted by gentle chemical lysis. Crude DNA extract was separated by gravity-flow gel filtration using a borosilicate glass column (BIORAD #737-0717). The column (7x150 mm) was packed with 1/6 volume PVPP (SIGMA # P-6755) and an upper layer of Sepharose 2B (SIGMA # 2B-300) to a final volume of 5 ml. Equilibration and removal of storage buffer was done by passing 5 reservoir volumes of TEN buffer (100 mM NaCl, 10 mM Tris, 1 mM EDTA; pH 8.0) by

gravity flow. Separation of DNA and humic acids was initiated by applying crude DNA extract (ideally 1-5 % resin volume) to the column. Approximately 300 µl fractions were collected dropwise (6 drops per fraction) and subsequently analysed. Relative DNA contents was quantified by fluorescence intensity comparison of DNA bands of each fraction on 1% agarose. Absorption at 230 nm and 260 nm was measured spectrophotometrically. Fractions 7 -10 represent pure DNA. Analysis of later fractions 11 – 22 show humic acid content with a maximum in fraction 16. As humic substances also absorb at longer wavelengths the OD<sub>260</sub> values increase in a similar extent as those of OD<sub>230</sub> which could mislead to assume high DNA concentrations. The specialist is able to discriminate pure DNA fractions from those contaminated by humic acids by increasing absorption at 230 nm and a low ratio OD<sub>260</sub>/OD<sub>230</sub>.

#### **Example 7: Validation of the purifying effect of the 2-phase gel on high molecular weight soil DNA**

Pellets of uncultivated bacteria isolated from soil were embedded in agarose plugs and treated appropriately to liberate their genomic DNA in situ. The DNA plugs were embedded either in an ordinary preparative PFGE (pulsed field gel electrophoresis) gel (Fig. 14 A) or a 2-phase PFGE gel (Fig. 14 C) for sizing and eventual purification. After completion of the PFGE, agarose plugs containing high molecular weight DNA in the compression zones were cut out and subjected to digestion with the restriction endonuclease *Hind* III to create DNA ends compatible with common cloning sites in suitable DNA vectors. After terminating the digestion the DNA plugs were inserted into analytical PFGE gels to analyse the success of the digestion.

As seen in Fig. 14 B soil DNA separated in a conventional agarose PFGE gel was highly resistant to digestion even with large amounts of restriction enzyme whereas high molecular weight DNA purified in the 2-phase PFGE gel was readily fragmented with very little restriction enzyme (Fig. 14D) proving its high degree of purity and consequently clonability. Cloning DNA fragments with compatible ends into suitably prepared DNA vectors like BAC, Plasmid, Lambda, Fosmid, Cosmid, YAC and PAC are standard procedures known to those skilled in the art.

Table 1:

Predicted RNA and protein encoding genes in the archaeal fosmid 29i4.

ORF	Nt range	protein size	putative function	most similar ortholog*	Phyl. aff.*	Comments
01	1-2367	789 aa (C-term. truncated)	Family B DNA polymerase pfam00136	AAC62689 Cenarchaeum symbiosum (0.0)	AEB	
RNA	3548-5018		16S RNA	85% identity to 16S RNA of Cenarchaeum symbiosum		
RNA	5847-8818		23S RNA	77% identity to 23S RNA of Cenarchaeum symbiosum		
02	8968-9849	293 aa	$\alpha/\beta$ hydrolase pfam00561	AAD02150 Pseudomonas stutzeri (e -17)	AEB	Catalytic triade Ser/Asp/His; close homologs are from bacteria
03	9888-10205	105 aa	Hypothetical	none		
04	10219-11304	361 aa	PHA Synthase ( $\alpha/\beta$ hydrolase, pfam00561)	P45366 Thiocystis violacea (e -66)	B	No homolog in archaea but PHB production has been described in Halobacteriaceae
05	11285-11866	193 aa	Hypothetical	none		
RNA	12414-12502		tRNA <sup>Asp</sup>			
06	13017-13454	145 aa	Hypothetical	none		
07	14324-15238	304 aa	Glycosyl transferases group 1, pfam00534	AAK41834 Sulfolobus solfataricus (e -11)	AEB	Transfer of ADP, UDP, GDP, CMP linked sugars
08	15716-17407	563 aa	Asparagine synthetase pfam 00310	AAB99117 Methanococcus Janaschii (e -67)	AEB	
09	17492-18157	221 aa	Phosphoserin phosphatase pfam00702	AAB86099 Methanothermobacter thermautotrophicus (e -14)	AEB	
10	18377-19588	403 aa	Conserved hypothetical	CAB50138 Pyrococcus abyssi (e -91)	A	Homologs only found in P. abyssi, P. horikoshii, A. permix, Glycosyltransferase group 2 domain; pfam 00535
11	20630-21793	387 aa	Transmembrane protein	BAB50489 Mesorhizobium loti (e -31)	B	Domain pfam 01173
12	24580-25461	293 aa	Fix A pfam01012	P53576 Azotobacter vinelandii (e -23)	AB	Paralogs of ETF $\beta$ : electron transfer flavoprotein $\beta$ subunit
13	25458-26741	427 aa	Fix B pfam00766	P53578 Clostridium saccharobutylicum (e -45)	AB	Paralogs of ETF $\alpha$ : electron transfer flavoprotein $\alpha$ subunit
14	26738-28615	625 aa	Fix CX	NP_454687 Salmonella enterica (e -36) FixX: Thermoplasma volcanium (e -08)	AB	Fused protein of FixC and FixX
15	29228-31465	745 aa	Sensory transduction histidine kinase, pfam00512	BAB73503 Nostoc sp. PCC 7120 (e -12)	AEB	
16	32505-33023	172 aa	Hypothetical protein	none		
17	32918-33925	335 aa (truncated)	hypothetical	none		

\*proteins are designated by their gene identification numbers followed by the species name. The e-values of blastx searches are added in brackets.

# Phyl. aff. = Phylogenetic affinity, denotes occurrence of homologs in Archaea (A), Bacteria (B) or Eucarya (E).